

Genetic barriers to historical gene flow between cryptic species of alpine bumblebees revealed by comparative population genomics

Matthew J. Christmas¹, Julia C. Jones^{1,2}, Anna Olsson¹, Ola Wallerman¹, Ignas Bunikis³, Marcin Kierczak⁴, Valentina Peona⁵, Kaitlyn M. Whitley^{6,7}, Tuuli Larva¹, Alexander Suh^{5,8}, Nicole E. Miller-Struttmann⁹, Jennifer C. Geib⁶, Matthew T. Webster¹

1) Department of Medical Biochemistry and Microbiology, Science for Life Laboratory, Uppsala University, Uppsala, Sweden

2) School of Biology and Environmental Science, University College Dublin, Dublin, Ireland

3) Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala University, Uppsala, Sweden

4) Dept of Cell and Molecular Biology, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Uppsala University, Uppsala, Sweden

5) Department of Organismal Biology – Systematic Biology, Uppsala University, Uppsala, Sweden

6) Department of Biology, Appalachian State University, Boone, North Carolina, USA

7) U.S. Department of Agriculture, Agriculture Research Service, Charleston, South Carolina, USA

8) School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich, UK

9) Biological Sciences Department, Webster University, St. Louis, Missouri, USA

Corresponding author: Matthew Webster, matthew.webster@imbim.uu.se

Abstract

Evidence is accumulating that gene flow commonly occurs between recently-diverged species, despite the existence of barriers to gene flow in their genomes. However, we still know little about what regions of the genome become barriers to gene flow and how such barriers form. Here we compare genetic differentiation across the genomes of bumblebee species living in sympatry and allopatry to reveal the potential impact of gene flow during species divergence and uncover genetic barrier loci. We first compared the genomes of the alpine bumblebee *Bombus sylvicola* and a previously unidentified sister species living in sympatry in the Rocky Mountains, revealing prominent islands of elevated genetic divergence in the genome that co-localize with centromeres and regions of low recombination. This same pattern is observed between the genomes of another pair of closely-related species living in allopatry (*B. bifarius* and *B. vancouverensis*). Strikingly however, the genomic islands exhibit significantly elevated absolute divergence (d_{xy}) in the sympatric, but not the allopatric, comparison indicating that they contain loci that have acted as barriers to historical gene flow in sympatry. Our results suggest that intrinsic barriers to gene flow between species may often accumulate in regions of low recombination and near centromeres through processes such as genetic hitchhiking, and that divergence in these regions is accentuated in the presence of gene flow.

Introduction

Genome-wide comparisons of genetic variation between species provide information about their history of divergence from a common ancestor. As populations diverge, barriers to gene flow eventually arise at multiple loci in their genomes (termed 'barrier loci'), which contain variants that govern ecological specialization or generate intrinsic genomic incompatibilities (Ravinet et al. 2017). Such barriers to gene flow may accumulate while gene flow is ongoing, such as in the case of sympatric or parapatric speciation, or alternatively in the absence of gene flow according to a strict allopatric model (Coyne and Orr 2004). Periods of gene flow can also occur when there is secondary contact between diverging species, during which barriers to introgression may either accumulate or break down (Kirkpatrick and Ravigné 2002; Rundle and Nosil 2005). When species hybridize, selection is predicted to act against gene flow at barrier loci but not in the rest of the genome (Wu 2001). However, despite intense study of many systems, we still lack a general understanding of which genomic regions tend to harbour barrier loci, how such barriers accumulate, and how the transition from incomplete to complete reproductive isolation occurs.

Comparisons of the genomes of closely-related species often reveal a heterogeneous landscape of divergence, which contain distinct peaks that have been described as islands of divergence (IoDs)

(Turner et al. 2005). This pattern has been interpreted according to several models. Firstly, if gene flow has been common between the species either during initial divergence in sympatry or after secondary contact, then IoDs could represent barrier loci where introgression is disadvantageous and selected against, leading to increased levels of divergence (Wu 2001). A large number of studies have used comparisons of genome-wide variation in recently-diverged species in order to identify IoDs, often with the aim of revealing genes that promote local adaptation and/or speciation and are recalcitrant to gene flow according to this model (Turner et al. 2005; Ellegren et al. 2012; Martin et al. 2013; Renaut et al. 2013; Poelstra et al. 2014; Soria-Carrasco et al. 2014; Lamichhaney et al. 2015; Malinsky et al. 2015; Chapman et al. 2016; Talla et al. 2017; Irwin et al. 2018; Papadopulos et al. 2019; Stankowski et al. 2019).

Secondly, in some cases it has been shown that IoDs represent ancient balanced polymorphisms that segregated in the ancestral populations (Guerrero and Hahn 2017; Han et al. 2017). Islands of divergence formed by this model are also expected to contain loci of adaptive significance. Such loci evolve under balancing selection in the ancestral population followed by sorting of divergent ancient haplotypes in the descendent populations.

Thirdly, IoDs with elevated relative divergence can form in the absence of gene flow via linked selection due to genetic hitchhiking or background selection, which has a greater effect in regions of reduced recombination (Nordborg et al. 1996; Charlesworth et al. 1997; Turner and Hahn 2010; Cruickshank and Hahn 2014). This process results in elevated levels of relative divergence (F_{ST}) in regions of low recombination between recently-diverged species that have not experienced gene flow. This can result in the formation of IoDs in regions of low recombination, termed "incidental islands", which do not harbour barrier loci or have a function in adaptation or speciation.

In some species comparisons, IoDs have been identified that are clearly associated with ecological specialisations, such as beaks of Darwin's finches (Lamichhaney et al. 2015; Han et al. 2017), or known incompatibilities, such as generation of melanoma in hybrids of swordtail fish (Powell et al. 2020). However, many studies have identified landscapes of divergence more consistent with the incidental island model, in which IoDs tend to occur in regions of low recombination (Turner and Hahn 2010; Cruickshank and Hahn 2014; Burri et al. 2015; Feulner et al. 2015; Ravinet et al. 2017; Talla et al. 2017) and may not be relevant for adaptation or speciation. It is however important to note that these processes are not mutually exclusive: IoDs formed due to linked selection or balancing selection could also harbour barrier loci (Ravinet et al. 2017) and the landscape of divergence could be shaped by multiple interacting processes (Chapman et al. 2016; Papadopulos et

al. 2019). Furthermore, there is also evidence from species that are known to hybridize that the rate of gene flow is correlated with recombination rate, which suggests that identification of IoDs in regions of low recombination does not preclude them from containing barriers to introgression (Schumer et al. 2018; Martin et al. 2019).

A key difference predicted between IoDs formed under these models is that those that have acted as barriers to gene flow or that are involved in ancient balanced polymorphisms should have longer coalescence times than the rest of the genome, resulting in IoDs with elevated absolute divergence measured by d_{XY} . Conversely, IoDs formed by ongoing linked selection in the absence of gene flow should not have elevated d_{XY} and may even have shorter coalescence times than the rest of the genome, due to these processes also occurring in the ancestral population, which results in reduced d_{XY} (Cruickshank and Hahn 2014; Irwin et al. 2018). Hence, the genomic landscape of absolute divergence reflects the presence or absence of historical gene flow or ancient balanced polymorphisms, and the existence of barrier loci. This landscape can be interrogated to learn about the processes that have occurred as species diverged.

Previous studies have demonstrated that regions of low recombination, particularly in the vicinity of centromeres, tend to accumulate elevated relative divergence in the absence of gene flow, due to the effects of linked selection (Nachman and Payseur 2012; Roesti et al. 2012). However, it is unclear whether this process also leads to the accumulation of barriers to gene flow in these regions, which could be important in establishing reproductive isolation. In addition, some models of divergence with gene flow where selection acts on many loci have also been shown to result in elevated divergence in regions of low recombination (Michel et al. 2010). Comparisons of patterns of divergence in pairs of species that have experienced different degrees of gene flow during their divergence can help reveal whether divergence in regions of low recombination is associated with the presence of barriers to gene flow.

Here we use population-scale genome sequencing to infer the mechanisms behind species divergence within the bumblebee subgenus *Pyrobombus* (Hines et al. 2006; Cameron et al. 2007; Martinet et al. 2019). Convergence in coloration due to Müllerian mimicry results in highly similar morphologies among bumblebee species (Williams 2007; Ezray et al. 2019), which rely mainly on chemical signalling for mate recognition (Goulson 2003). These factors make species recognition particularly difficult in bumblebees and studies utilizing multiple genetic loci have recently resulted in the discovery of several previously undescribed species (Martinet et al. 2019; Ghisbain et al. 2020). The number of bumblebee species is likely underestimated with many cryptic species living in

sympatry, which may have experienced gene flow during their formation (Bertsch et al. 2004; Murray et al. 2008; Bossert 2015).

We constructed a highly-contiguous genome assembly of the species *Bombus sylvicola*, and surveyed genomic variation in this species by whole-genome resequencing 284 samples from across the Rocky Mountains in Colorado, USA. Unexpectedly, these samples fell into two distinct genetic clusters, revealing the presence of a previously unknown cryptic species living in sympatry with *B. sylvicola*, which we name *B. incognitus*. We performed genome-wide comparisons between these two sympatric species and contrasted them with genomic divergence between another pair of closely-related species living mainly in allopatry (*B. bifarius* and *B. vancouverensis*) to uncover evidence for historical gene flow and identify and characterise regions of the genome that have likely acted as barriers to gene flow in the past. Analysis of the genomic landscape of divergence reveals signals of gene flow in the sympatric but not the allopatric pair, and provides important insights into how genomic architecture influences the formation of barrier loci.

Results

A highly-contiguous genome assembly of Bombus sylvicola

We used a combination of Oxford Nanopore (ONT) and 10x Chromium sequencing to generate a genome assembly of the bumblebee *B. sylvicola* using a single haploid drone sample for each technology (see Methods). A recent study analysed genetic and morphological differentiation between *B. sylvicola* collected in northern Alaska and *B. lapponicus* collected in northern Sweden (Martinet et al. 2019). Based on relatively low levels of divergence, this study redefined *B. sylvicola* as the subspecies *B. lapponicus sylvicola*. However, here we maintain the previous name *B. sylvicola* for our samples for consistency with previous ecological studies in this region and because their relationship to the populations in Alaska has not been directly tested. The sequencing and assembly resulted in a set of 592 contigs with a total length of 252,081,862 bp and an N50 of 3,020,754 bp. Analysis of genome completeness estimated that 97.9% (5,865) of BUSCO genes were complete in the assembly and only 1.5% (92 genes) were undetected. We estimated the position and orientation of contigs on chromosomes via a whole-genome alignment with the *B. terrestris* genome. This resulted in the preliminary placement of 91.1% of the *B. sylvicola* genome onto 18 likely chromosomes (hereafter ‘pseudochromosomes’). Our annotation pipeline annotated 11,585 genes across the *B. sylvicola* genome (14% of the genome is found in exons). This is highly comparable to the 11,874 genes in the *B. terrestris* gene set (v. 1.0) and the 12,728 genes in the *B. impatiens* gene set (v. 2.1).

We analysed both the raw ONT reads and the assembled contigs to identify putative centromere-associated sequences (see Methods). This analysis revealed the presence of three 15 bp monomers at high frequency in the ONT reads, each differing by one to two base pairs and forming tandem repeat arrays (Supplementary Figure S1). Repeat arrays occur near or at contig boundaries, indicating that the true arrays are likely longer and that the assembler has failed to assemble them further due to their repeat nature (Kolmogorov et al. 2019). The prevalence and location in the vicinity of assembly gaps suggest that the repeats occur near centromeres. We find 83 occurrences of the tandem repeat array in lengths ranging from 28 bp to 22,416 bp (mean length = 1,911 bp) across the genome, with 48 occurrences across 15 of the 18 pseudochromosomes (the remaining 35 are on unplaced contigs). Nine of the 18 pseudochromosomes contain arrays greater than 1 kbp in length and, in most cases, arrays occur in one region per pseudochromosome, indicating the likely locations of centromeres.

Population-scale sequencing leads to identification of a new species

We collected 284 female worker bees identified phenotypically as *B. sylvicola* and 17 identified as *B. bifarius* from seven localities in the Rocky Mountains, Colorado, USA (Figure 1A). We obtained Illumina whole genome sequencing (WGS) data for all samples. We also obtained published WGS data from 4 samples of *B. bifarius* and 17 samples of *B. vancouverensis* collected from Colorado across north-eastern USA (Ghisbain et al. 2020) and 21 samples of *B. melanopygus* from western USA (Tian et al. 2019) giving a total of 343 re-sequenced genomes of bumblebees within the *Pyrobombus* subgenus (Figure 1A). We mapped these WGS datasets to our *B. sylvicola* genome assembly and performed variant calling. The mean coverage across all samples was 14.7x and we inferred 15,094,475 SNPs (see Supplementary Tables S1 and S2 for full details of all samples).

A principal-components analysis (PCA) of the genome-wide SNP dataset showed clear clustering by species (Figure 1B). Surprisingly, the 284 samples identified as *B. sylvicola* were split into two distinct clusters, containing 217 and 67 samples respectively, with no observations of intermediates between the two clusters. The *B. bifarius* and *B. vancouverensis* samples also formed two distinct clusters, consistent with their assignment as two separate species by Ghisbain et al (2020). A neighbour-joining tree also strongly supported the division of the *B. sylvicola* samples into two clusters with the *B. bifarius* – *B. vancouverensis* pair placed distantly from these clusters (Figure 1C). We also generated a neighbor-net network based on SNPs across the genome to check for any conflicting signals or alternative phylogenetic histories (Supplementary Figure S2), which demonstrates that the

underlying evolutionary history of these species is treelike. Taken together, these data indicate the presence of a cryptic species within the purported *B. sylvicola* samples.

We next attempted to reveal the identity of this cryptic species. *Bombus melanopygus* appears as an outgroup to the two *B. sylvicola*-like clusters, placing the cryptic species within the *B. lapponicus* - *B. sylvicola* - *B. monticola* species complex (Cameron et al. 2007; Williams et al. 2014; Martinet et al. 2019). Martinet et al. (2019) used both morphological and genetic analysis to delineate relationships of all species within this complex and identified a previously-undescribed cryptic species with similar morphology to *B. sylvicola*, which they call *B. interacti*, which could potentially match the unexpected cluster we identified. In order to test this possibility, we compared sequences of our samples at the PEPCK and COI loci with those of all species in the complex presented in Martinet *et al.* (2019).

We determined the PEPCK sequences of all of our samples from the whole-genome sequencing data (Supplementary Figure S3A) and found that samples from our larger *B. sylvicola*-like cluster closely corresponded to the *B. sylvicola* and *B. lapponicus* samples from Martinet et al. (2019), with three samples matching exactly. However, our smaller *B. sylvicola*-like cluster did not show similarity with any other species examined. It is most similar to *B. sylvicola* and clearly distinct from *B. interacti*. We next generated COI sequences from a subset of our samples using PCR and Sanger sequencing (Supplementary Figure S3B). Here we find a similar pattern, where the samples from our larger *B. sylvicola*-like cluster again correspond most closely with the *B. sylvicola* and *B. lapponicus* samples from Martinet et al. (2019), whereas samples from the smaller *B. sylvicola*-like cluster are closely related but distinct from these samples, and the *B. interacti* samples are more distant. Taken together, these results suggest that the samples from the larger *B. sylvicola*-like cluster in our dataset correspond to samples identified as *B. sylvicola* in previous studies, whereas the samples in the smaller *B. sylvicola*-like cluster represent a previously undescribed species. We assign this species the provisional name *Bombus incognitus*.

We performed a detailed characterisation of the anatomical structures on the heads and abdomens of a subset of samples identified genetically as *B. sylvicola* and *B. incognitus* (see Methods; Supplementary Table S3). None of these traits could be used to distinguish between the two species. *Bombus sylvicola* samples were significantly larger on average based on measurements of intertegular distance, which is a proxy of body size (mean = 3.84 mm and 3.59 mm for *B. sylvicola* and *B. incognitus* respectively; Wilcoxon rank sum test, $W = 7673$, $p = 1.304 \times 10^{-5}$; Supplementary Figure S4). The locations where samples of each species were found showed substantial overlap (Supplementary Table S2). Both species were collected together on six out of the seven localities (*B.*

incognitus was not present among the 28 samples collected on Mount Democrat). The two species were found at overlapping elevations, although there was a significant tendency for *B. incognitus* to be found at lower elevations (mean elevations 3,792 m and 3,664 m for *B. sylvicola* and *B. incognitus* respectively; Wilcoxon rank sum test, $W = 11192$, $p = 7.05 \times 10^{-12}$, Supplementary Figure S4). In summary, *B. sylvicola* and *B. incognitus* are cryptic species found in sympatry across our sampling localities.

The genomic landscape of divergence differs between sympatric and allopatric species pairs

We used our dataset to compare genome divergence on multiple spatial and temporal scales. Firstly, the two most geographically and genetically distant (as revealed by a PCA of all *B. sylvicola* samples; Supplementary Figure S5) populations of *B. sylvicola* that we sampled were from Niwot Ridge and Quail Mountain, Colorado (hereafter “within-species pair”). Secondly, we sampled two closely-related species existing in sympatry across this range (*B. sylvicola* and *B. incognitus*; hereafter “sympatric pair”), which allowed us to investigate genome divergence in species that may have had the opportunity to undergo gene flow in the past. Thirdly, we used genomic data for two other closely-related Pyrobombus species that exist mainly in allopatry, (*B. bifarius* and *B. vancouverensis*; hereafter “allopatric pair”), which we hypothesize underwent speciation in the absence of gene flow. These species are separated geographically with only a narrow range of overlap in mountains in Utah, USA, with no evidence of hybridisation (Ghisbain et al. 2020).

In order to compare landscapes of divergence at these different scales we carried out genome-wide sliding-window F_{ST} scans. Average genome-wide F_{ST} was 0.02, 0.41, and 0.14 for the within-species, sympatric, and allopatric pairs respectively. Both the within-species and allopatric pairs displayed typical F_{ST} distributions of a single large peak centred close to the median score with a tail representing relatively few regions with heightened divergence (Figure 2A, 2C). However, the sympatric pair displayed a striking bimodal distribution of F_{ST} , with a large peak centred at 0.21 followed by a second peak of extreme divergence centred at 0.93 (Figure 2B). The sympatric pair therefore stands out as having a distinct portion of the genome with highly elevated divergence. This pattern has been observed for other pairs of species that diverged under conditions of gene flow (Seehausen et al. 2014).

We used estimates of genetic variation and divergence to estimate the effective population size and timings of the splits between species (Table 1). *Bombus incognitus* showed the highest levels of

genetic variation despite being less abundant than *B. sylvicola* in our sampling localities, and both of these species exhibited higher levels of variation than *B. bifarius* and *B. vancouverensis*. Estimates of N_e are in the same range as estimated for honeybees (Wallberg et al. 2014). Average F_{ST} between *B. sylvicola* and *B. incognitus* based on all SNPs located outside of regions of extreme divergence was 0.34. This translates into an estimated divergence time of $t = 396,000$ (95% c.i. 389,000 – 403,000) generations since the species split under a simple demographic model using estimates for θ_w per base and N_e (Table 1). For the allopatric pair, average F_{ST} is 0.12, which indicates a divergence time of $t = 67,290$ (95% c.i. 66,177 – 68,514) generations. Given the generation time of bumblebees is one generation per year, these divergence times translate directly into years.

Genomic islands of divergence evolve in similar locations in independent species comparisons

We converted F_{ST} values to Z-scores (ZF_{ST} ; Figure 3A-C) and used them to define ‘highly divergent windows’ in each comparison, where $ZF_{ST} \geq 2$ (2 or more SDs above the median F_{ST}). This resulted in 486 highly divergent windows (3.86% of the genome) for the within-species pair, 1,758 windows (13.95% of the genome) for the sympatric pair, and 842 windows (6.68% of the genome) for the allopatric pair. It is strikingly apparent from the genome-wide ZF_{ST} plots (Figure 3A-C), particularly for the sympatric pair (Figure 3B), that genome variation contains several large blocks of extreme divergence. We defined blocks larger than 100 kbp as islands of divergence (IoDs; see Methods). Using this definition there are 20 IoDs in the within-species pair (Figure 3A; average length of 560 kbp), 28 IoDs in the sympatric pair (Figure 3B; average length of 1.26 Mbp) and 68 IoDs in the allopatric pair (Figure 3C; average length of 223 kbp). The longest 18 IoDs in the sympatric pair were positioned over 17 pseudochromosomes, demonstrating that, in the majority of cases, pseudochromosomes contain a single major IoD in sympatry (Figure 3B). A similar pattern is observed in the within-species and allopatric pairs, although IoDs are smaller and not found on all pseudochromosomes. The sympatric pair is therefore distinguished by the presence of one large IoD per pseudochromosome and having the greatest proportion of the genome in IoDs.

There was a highly significant overlap in the location of IoDs between the independent population/species comparisons, assessed using permutation tests (Figure 3D; Supplementary Figure S6). Notably, 75% of within-species IoDs overlapped with a sympatric IoD (permutation test, Z-score = 5.194, $p = 0.001$; Supplementary Figure S6A) and 43% of allopatric IoDs overlapped with a sympatric IoD (permutation test, Z-score = 4.567, $p = 0.001$; Supplementary Figure S6B). The within-species and allopatric comparisons showed some overlap but it was not significant at $p=0.05$

(permutation test, Z-score = 1.576, $p=0.09$; Supplementary Figure S6C). The observation of IoDs in the within-species comparison shows that divergence can accumulate in these regions in the absence of reproductive isolation.

Islands of divergence are associated with reduced recombination rate and a common satellite repeat

We estimated genome-wide variation in recombination rate in *B. sylvicola* using LDhat (Figure 3E) (McVean et al. 2004). We found a significant positive correlation between recombination rate and GC content (Pearson's $r = 0.23$, $p < 0.001$) which has been found previously in bumblebees (Kawakami et al. 2019) and is observed widely in sexual eukaryotes (likely as a result of GC-biased gene conversion in high-recombining regions) (Pessia et al. 2012). We found a strongly significant negative correlation between ZF_{ST} and recombination rate that was particularly pronounced in the sympatric pair (Spearman's $\rho = -0.78$, $p < 2.2 \times 10^{-16}$) but lower in the within-species and allopatric pairs (within-species: Spearman's $\rho = -0.02$, $p = 0.03$; allopatric: Spearman's $\rho = -0.24$, $p < 2.2 \times 10^{-16}$). In addition, we found significantly reduced recombination rates in IoDs in all species pairs (Wilcoxon rank sum test, $p < 2.2 \times 10^{-16}$ in all cases; Supplementary Figures S7, S8). This finding is less pronounced in the allopatric pair, where regions of high divergence are found to occur in regions with higher recombination rates. We note however that as the method we employed estimates recombination rates based on patterns of genetic variation it is likely that the extremely low levels of nucleotide diversity in the IoDs in *B. sylvicola* result in lower accuracy to measure recombination rates in these regions.

In all three comparisons we found significantly lower GC content, lower mappability (a measure of sequence uniqueness in the genome) and higher repeat content inside IoDs compared to the rest of the genome (Wilcoxon rank sum tests, all significant at $p < 2.2 \times 10^{-16}$; Table 2; Supplementary Figure S7). These observations are consistent with a tendency for IoDs to occur in regions of low recombination. We also observed significantly greater gene density inside IoDs compared to random expectations in the allopatric comparison (Table 2). In the sympatric comparison, IoDs comprise 13.9% of the genome and contain 2,135 genes (19.8%). It is plausible that the high levels of divergence across this large number of genes in IoDs have functional consequences that could contribute to adaptation or intrinsic barriers to gene flow.

We next tested for associations between the locations of the likely centromeric tandem repeat arrays and the locations of IoDs in each population/species comparison (Figure 3F; Supplementary

Figure S9). There is a particularly strong association in the sympatric pair, where we observe one dominant large loD (>100 kbp) per pseudochromosome. We found a significant overlap of 36 (75%) repeats overlapping with 13 sympatric loDs, which was 4.6x greater than expected by chance (permutation test, Z-score = 6.33, $p = 0.001$). Significant overlap was also observed in the within-species comparison, where 12 (25%) repeats overlapped with four loDs (permutation test, Z-score = 3.25, $p = 0.02$). In the allopatric pair, there was overlap between loDs and repeats with 8 (17%) repeats overlapping 6 loDs, but this was not significantly different from the overlap expected by chance (permutation test, Z-score = 2.16, $p = 0.06$). Hence, although there is a tendency for loDs to occur near centromeres in all comparisons, there is a particularly strong association in the sympatric comparison. It is unlikely that errors in read mapping or variant calling in repetitive regions contributed to the elevated divergence inferred in loDs. Inspection of loDs reveals that elevated divergence occurs both in repetitive regions and in flanking non-repetitive regions (Supplementary figure S10).

In order to test whether loDs may represent large structural inversions, we ran the program manta (Chen et al. 2016) on a subset of the *B. sylvicola* and *B. incognitus* bam files. We found evidence for only three short inversions that were fixed between the two species: one of 6,351 bp on contig_013, one of 3,059 bp on contig_026, and one of 1,691 bp on contig_118. All three of these putative inversions are found within loDs and likely lead to some of the high divergence we see in these regions, however they make up only a small fraction of the loDs they are found in. We therefore did not find support for any of the identified loDs representing structural rearrangements.

Genomic islands of divergence have elevated d_{XY} in the sympatric but not the allopatric pair

For all species/population comparisons, nucleotide diversity (π) was significantly lower in loDs compared to the rest of the genome (Wilcoxon rank sum test, $p < 2.2 \times 10^{-16}$ in all cases; Figure 4A-C; Supplementary Figures S11-13). This is consistent with the action of linked selection (background selection and/or genetic hitchhiking) on these regions, which both increases relative divergence and decreases levels of genetic variation. There is a more pronounced reduction of π in loDs in both sympatric species (80% and 77% decrease in π inside compared to outside loDs in *B. sylvicola* and *B. incognitus* respectively) compared to the allopatric species (28% and 30% decrease in π inside compared to outside loDs in *B. bifarius* and *B. vancouverensis* respectively). loDs are therefore more extensive and exhibit lower within-species variation in the sympatric comparison.

We found d_{XY} to be significantly lower in IoDs compared to the rest of the genome in the within-species and allopatric pairs (Wilcoxon rank sum tests, $p < 2.2 \times 10^{-16}$ in both cases), with a 66% and 48% reduction in d_{XY} inside IoDs respectively. Yet intriguingly, we found d_{XY} to be significantly elevated inside compared to outside IoDs in the sympatric pair (Wilcoxon rank sum test, $p < 2.2 \times 10^{-16}$), with a 17% increase inside IoDs (Figure 4D-F). Neighbour-joining trees based on only SNPs found a) within IoDs and b) outside of IoDs revealed the same topology, but substantially longer branch lengths in the sympatric pair within IoDs, consistent with their elevated d_{XY} (Figure 4G,H). The sympatric comparison is therefore distinguished by a strikingly bimodal distribution of F_{ST} across the genome, and IoDs with significantly elevated d_{XY} .

Two main evolutionary scenarios have been demonstrated to result in elevated d_{XY} in IoDs. The first is divergence with differential gene flow (Cruickshank and Hahn 2014). The second is the presence of ancient balanced polymorphisms in the ancestral population that are sorted in the descendent populations (Guerrero and Hahn 2017). However, the association of IoDs with elevated d_{XY} in the sympatric but not the allopatric comparison is most parsimoniously explained by differences in the incidence of gene flow between the two comparisons, suggesting that the sympatric IoDs have been shaped by differential gene flow (see Discussion). The observation of reduced d_{XY} in IoDs in the allopatric pair is most consistent with ongoing linked selection resulting in "incidental islands".

Differences in mutation rate or average levels of evolutionary constraint in IoDs compared to the rest of the genome could also potentially influence the differences in d_{XY} we observed in these regions. Such differences would generate variation in d_{XY} in more distant species comparisons. In order to assess this possibility, we estimated d_{XY} in the IoD regions defined by the sympatric comparison in comparisons of both *B. sylvicola* and *B. incognitus* to *B. bifarius* and to a single sample of *B. balteatus*, which belongs to a separate subgenus. In all comparisons, d_{XY} inside and outside of IoDs did not significantly differ (Wilcoxon rank sum test, $p > 0.05$; supplementary figure S14). This indicates that the average rate of nucleotide substitution in IoDs is similar to the rest of the genome.

Islands of divergence are associated with extended drops in levels of genetic variation in the sympatric comparison

We calculated the population branch statistic (PBS) (Yi et al. 2010) to assess the relative amount of divergence that had occurred along each branch. Population branch statistic correlated strongly between branches in the sympatric pair, (Spearman's $\rho = 0.62$, $p < 2.2 \times 10^{-16}$) as well as with F_{ST} in both comparisons (Spearman's $\rho = 0.85$ and 0.60 in the sympatric and allopatric comparisons

respectively; $p < 2.2 \times 10^{-16}$ for in both comparisons). As expected, IoDs have significantly greater PBS compared to the rest of the genome in all four species (Wilcoxon rank sum test, $p < 2.2 \times 10^{-16}$ in all cases), showing there to be strong correspondence between the locations of regions of elevated divergence that have formed on both branches leading from the common ancestor in both species pairs (Supplementary Figure S15).

To further characterize IoDs between our pairs of species, we calculated average π , ZF_{ST} , and d_{XY} with increasing distance from their midpoints (Figure 5). There is a more extensive reduction of π in IoDs in the sympatric pair compared to the allopatric pair (Figure 5A). For the sympatric species, average π values remained below the genome average up to ~ 1.5 Mbp away from the centre of the IoDs, however for the allopatric species this distance was ~ 0.2 Mbp. A more extensive drop in within-species variation at IoDs in the sympatric pair is unlikely to be explained by a greater effect of linked selection and suggests the influence of differential gene flow.

Average ZF_{ST} values reflected the larger size of IoDs in the sympatric pair compared to the within-species and allopatric pairs: in the sympatric pair, ZF_{ST} did not return to the genome average until ~ 1.8 Mbp away from the centre of the IoDs, whereas for the within-species pair and allopatric pair the distance was ~ 1 Mbp (Figure 5B). There is a stark contrast in patterns of d_{XY} in IoDs in the sympatric compared to the within-species and allopatric pairs (Figure 5C). For the sympatric pair, average d_{XY} values remained *above* the genome average up to ~ 2.6 Mbp away from IoD centres, whereas for the within-species and allopatric pairs d_{XY} values remained *below* the genome average until ~ 1.8 Mbp and 0.4 Mbp away from IoD centres respectively. IoDs in the sympatric comparison are therefore larger, show more extensive reductions in genetic variation, and are distinguished by elevated absolute divergence (d_{XY}). These observations are all consistent with a scenario where IoDs in the sympatric comparison have acted as barriers against historic gene flow. In contrast, the lower d_{XY} in IoDs in the within-species and allopatric pairs likely reflect the action of linked selection in the ancestral populations, which reduces coalescence times between species (Irwin et al. 2018).

Discussion

We analysed genome variation in multiple closely-related species of bumblebees in the *Pyrobombus* subgenus to uncover mechanisms of species divergence and isolation. Analysis of 284 specimens classified as *B. sylvicola* in the Rocky Mountains revealed a previously-undetected cryptic species living in sympatry, which we call *B. incognitus*. Genome-wide comparisons of genetic variation between *B. sylvicola* and *B. incognitus* revealed a striking bimodal landscape of divergence, with

extensive pronounced genomic islands of divergence (IoDs) in the vicinity of centromeres. Our analysis indicates that these centromere-associated IoDs contain barrier loci that have restricted gene flow in these regions between *B. sylvicola* and *B. incognitus* in the past, whereas gene flow was able to continue to a greater extent elsewhere in the genome. We find no evidence of contemporary hybridization between these species, suggesting that current gene flow is rare or non-existent. Thus, our findings provide a window into the processes that lead to reproductive isolation.

Bombus incognitus is a previously-undetected bumblebee species

The newly-discovered species *B. incognitus* was indistinguishable from *B. sylvicola* using typical diagnostic characters used to identify *B. sylvicola*, and all *B. incognitus* samples were initially classified as *B. sylvicola* on the basis of morphology. Analysis of the population genomic data indicated the presence of two distinct clusters with high divergence ($F_{ST} = 0.41$) across most sampling locations, strongly indicating the presence of two species. A more detailed comparison of the head anatomy (see Methods) did not reveal any characters that distinguish between these two species. However, we observed that the *B. incognitus* we collected were on average 6.6% smaller than *B. sylvicola*, and were found at lower altitudes despite an overlapping range.

Our whole-genome comparisons demonstrate that *B. sylvicola* and *B. incognitus* are more closely related to each other than either is to *B. melanopygus*. This places *B. incognitus* within a clade of the *Pyrobombus* subgenus that also contains *B. bimaculatus*, *B. monticola*, *B. konradini*, *B. lapponicus* and *B. interacti*. *Bombus lapponicus* is an Old-World species with very low divergence in morphology and genetic distance from the New World *B. sylvicola*. Due to this observation, Martinet et al. suggested that *B. sylvicola* should be considered a subspecies of *B. lapponicus* (*B. lapponicus sylvicola*) (Martinet et al. 2019). The species *B. interacti* shows a strong resemblance to *B. sylvicola*. We constructed phylogenies of sequences from the PEPCK and COI loci from our samples compared to samples of all species in the clade presented by (Martinet et al. 2019). Our analyses show a strong correspondence between our *B. sylvicola* samples from Colorado and samples of this species collected in Alaska. However, our *B. incognitus* samples do not cluster together with *B. interacti*, or match any other known species, and are most closely related to the *B. sylvicola* samples. This indicates that *B. incognitus* differs from previously-described members of this clade and should be considered a separate species, further adding to the amount of cryptic species diversity recognised in this clade. Further studies are necessary to comprehensively characterise the range and morphology of this species.

We considered the possibility that *B. incognitus* is a hybrid species resulting from interbreeding between *B. sylvicola* and *B. melanopygus* or another more distantly related species. Under this scenario, the majority of the genome, which exhibits low divergence between *B. sylvicola* and *B. incognitus*, would derive from their common ancestor whereas the highly-divergent IoDs would have a shorter evolutionary distance to *B. melanopygus* in *B. incognitus*. However, Figure 4 shows that this is not the case, as regions defined as IoDs have longer branches throughout the tree and relationships between species are the same within and outside of IoDs. Furthermore, examination of PBS across the genome (Supplementary Figure S10) indicates that the level of divergence is similar along the branches leading to *B. sylvicola* and *B. incognitus*, both within and outside of IoDs, indicating that these species diverged from a common ancestor and neither resulted from a hybridization event from a more distantly related species. Finally, there is no indication of hybridisation from the neighbor-net network constructed from the dataset (Supplementary Figure S2).

The lack of intermediates among samples of *B. sylvicola* and *B. incognitus* indicates that the species do not commonly hybridize, although patterns of genomic variation in both species indicate gene flow was ongoing for some time during their divergence. Extensive genome resequencing using field collections could lead to the discovery of more cryptic species in many taxa in future. Such discoveries are arguably more likely in bumblebees because their morphologies often converge due to Müllerian mimicry (Williams 2007; Ezray et al. 2019) and mate recognition occurs mainly via chemical signals (Goulson 2003), which could mask species diversity. Accurate identification of distinct clusters and divergence times requires data from multiple loci or the whole genome due to the effects of incomplete lineage sorting producing conflicting phylogenetic signals (Mallet et al. 2016).

The genomic landscape of divergence can be shaped by multiple factors

Comparisons of genome-wide variation in recently-diverged species often reveal a highly variable landscape of divergence containing IoDs (Turner et al. 2005; Ellegren et al. 2012; Renaut et al. 2013; Poelstra et al. 2014; Soria-Carrasco et al. 2014; Lamichhaney et al. 2015; Malinsky et al. 2015; Chapman et al. 2016; Talla et al. 2017; Irwin et al. 2018; Papadopulos et al. 2019; Stankowski et al. 2019; Liu et al. 2020) . One interpretation of IoDs is that they contain barrier loci that hinder introgression in species that have experienced gene flow (Wu 2001). Under this scenario, gene flow

is prevented at IoDs but continues in the rest of the genome, leading to elevated divergence at IoDs. Two main models of the accumulation of divergence under gene flow have been proposed, which can be viewed as extremes on a continuum (Feder et al. 2012). Under the classical island view, IoDs contain specific outlier loci that form barriers to gene flow. Elevated divergence spreads to nearby tightly-linked loci through divergence hitchhiking, whereby variants in linkage with these barrier loci are also prevented from introgressing (Via and West 2008; Via 2012).

Under the second view of divergence under gene flow, termed the continent view of genomic divergence, a much larger number of selected loci contribute to genetic isolation across the genome (Feder et al. 2012). The landscape of divergence is shaped by selection at these loci mediated by the genomic architecture, which includes factors such as linkage relationships, recombination rate variation and the strength of selection at each locus (Michel et al. 2010; Feder et al. 2012). This second view implies that species barriers are polygenic and does not rely on an effect of divergence hitchhiking. Islands of divergence (or "continents") generated by this model are larger and more likely to occur in regions of low recombination, such as inversions (Yeaman 2013) or at centromeres (Turner et al. 2005; Liu et al. 2020) because the effects of linked selection are enhanced in these regions.

It has also been demonstrated by both theoretical and empirical studies that IoDs can evolve in the absence of gene flow. One mechanism that can generate IoDs is sorting of ancient balanced polymorphisms, in which divergent haplotypes segregating in an ancestral population become fixed after the populations split (Guerrero and Hahn 2017). This mechanism likely occurred at loci that govern beak morphology in Darwin's finches, at which ancient haplotypes that predate the split between species are observed (Han et al. 2017). Another mechanism that can generate IoDs is linked selection (genetic hitchhiking and background selection) mediated by genomic architecture, which causes IoDs to recurrently arise in regions of low recombination (Cruickshank and Hahn 2014; Burri et al. 2015). This phenomenon is connected to the observation across taxa that genetic variation is reduced in regions of low recombination because of linked selection (Begun and Aquadro 1992). This effect causes measures of relative divergence such as F_{ST} to be elevated in regions of low recombination. IoDs generated by this effect do not necessarily contain loci involved in adaptation or barriers to gene flow. However, it is possible that IoDs generated by linked selection can harbour divergent loci that subsequently act as barriers under conditions of gene flow.

Genomic barriers to gene flow between B. sylvicola and B. incognitus

Here, we aimed to distinguish between these scenarios and identify putative barrier loci using independent comparisons of genomic variation: a) between populations of *B. sylvicola* from different localities, b) between *B. sylvicola* and *B. incognitus* populations living in sympatry, and c) between two additional species (*B. bifarius* and *B. vancouverensis*) living mainly in allopatry. We find that IoDs, defined by elevated F_{ST} , recur in the same genomic locations associated with low recombination and centromeres in all three independent comparisons. Strikingly however, we find that the genomic landscape of divergence in the sympatric comparison displays distinct features that indicate that it has been shaped by differential gene flow. Firstly, a measure of absolute divergence, d_{XY} , is elevated in IoDs in the sympatric but not the allopatric comparison. The d_{XY} statistic has been shown by simulations to be elevated in IoDs under conditions of differential gene flow (Cruickshank and Hahn 2014). Secondly, there is a more extensive drop in genetic variation in IoDs (~2 Mb on average) in both species of the sympatric pair compared to the allopatric pair. This may indicate a greater effect of selection in removing introgressed alleles at barrier loci in these regions. Thirdly, there is a markedly bimodal distribution of window-based F_{ST} in the sympatric pair, which is not present in the other comparisons, also reflecting the presence of extensive IoDs in this comparison. A similar distribution of F_{ST} has also been observed among *Heliconius* species that have undergone speciation without geographical isolation (Martin et al. 2013; Seehausen et al. 2014).

Another process that can result in IoDs with elevated d_{XY} is balancing selection in the ancestral population (Guerrero and Hahn 2017; Han et al. 2017). Although we cannot formally exclude the possibility that IoDs in the sympatric comparison represent ancient balanced polymorphisms, several observations strongly favour the differential gene flow model. Firstly, ancient balanced polymorphisms are not expected to be strongly associated with centromeres and regions of low recombination, which is observed for IoDs in all comparisons. Balanced polymorphisms are expected to be associated with loci involved in adaptation, which should not be strongly biased in their genomic locations. Secondly, we would not expect IoDs in the same genomic locations to represent ancient balanced polymorphisms in the sympatric comparison but not the allopatric one, considering that ancient balanced polymorphisms would be able to sort regardless of the presence of gene flow. Elevated d_{XY} in IoDs in the sympatric comparison therefore most likely indicates that they have acted as barriers to introgression during periods of gene flow.

The IoDs identified in the allopatric comparison appear most consistent with formation by linked selection in the absence of gene flow. No evidence for recent gene flow in the allopatric comparison is revealed by our data, as also found by (Ghisbain et al. 2020), although historical gene flow or low

levels of ongoing gene flow where the ranges of the species overlap cannot be ruled out. The mechanism by which IoDs formed in the sympatric comparison is less clear. One possibility is that they were also formed by linked selection during periods when the species were isolated from each other. Another possibility is provided by the "continents" model of divergence-with-gene-flow, whereby selection against introgression at a large number of barrier loci across the genome mediated by the recombination landscape in the face of gene flow leads to IoDs (or continents) in regions of low recombination (Michel et al. 2010). More detailed modelling could potentially determine which of these scenarios is more feasible.

Our divergence time estimate of ~396,000 years between *B. sylvicola* and *B. incognitus* coincides with a period of global cooling that was followed by rapid global warming around 340,000 years ago (Vimeux et al. 2002; Uemura et al. 2018). This could have been a driver of subpopulation isolation as cold-adapted alpine species likely became more isolated in mountain top habitats under warming (Hewitt 2000; Hines 2008). A subsequent period of cooling would then allow for secondary contact. It is therefore possible that a period of partial or complete geographic isolation facilitated the build-up of genetic incompatibilities by linked selection at IoDs which were then accentuated due to gene flow elsewhere in the genome during secondary contact. Quaternary climate oscillations are likely responsible for divergence in reproductive traits between populations of the red-tailed bumblebee (*B. lapidarius*) in Europe (Lecocq et al. 2013). Local fragmentation followed by gene flow during secondary contact could potentially be a common mode of speciation in high-altitude bumblebees, giving rise to cryptic species. Further modelling-based studies and empirical studies of more species are required to determine the validity and generality of this scenario. Knowledge of the distribution, ecology and population history of *B. incognitus* is currently completely lacking. However, while details of the speciation process are unclear, our evidence suggests that the evolution of barrier loci in extended regions of low recombination near centromeres has promoted reproductive isolation between these two species.

Recombination mediates the accumulation of barriers to gene flow

Our results are compatible with other studies demonstrating that hybridizing natural populations harbour numerous genetic incompatibilities throughout their genomes (Schumer et al. 2014). Reduced introgression in regions of low recombination has been observed in hybrids of swordtail fish (Schumer et al. 2018). Similarly, a correlation between recombination rate and introgression has also been inferred in *Heliconius* butterflies (Martin et al. 2019), *Mimulus* monkey-flowers (Brandvain et al.

2014), house mice (Janoušek et al. 2015), and between Humans and Neanderthals (Juric et al. 2016). These observations could be due to the interaction of linkage and selection against introgression of genetic incompatibilities (Coughlan and Matute 2020). Selection against genetic incompatibilities in regions of low recombination is expected to remove introgressed alleles in a larger portion of the genome due to linkage (Schumer et al. 2018). This mechanism could also promote differentiation in regions of low recombination under conditions of gene flow. Regions of low recombination could also accumulate fixed genetic incompatibilities in the absence of gene flow due to the effects of linked selection, which could also lead to reduced introgression in these regions upon secondary contact (Ravinet et al. 2017).

Bumblebees may be a particularly good model systems to uncover the influence of genome architecture on species divergence due to their extremely high rates of recombination. The average recombination rate in the bumblebee *B. terrestris* has been estimated as ~9 cM/Mb (Kawakami et al. 2019). High recombination rates have been estimated in other social insects, and appear to correlate with the degree of sociality (Wilfert et al. 2007) (rates in the highly social honeybee have been estimated as >20 cM/Mb (Kawakami et al. 2019)). Importantly, in both honeybees and bumblebees, there is also an extreme reduction in recombination rates in centromeres (Kawakami et al. 2019) and we observe a clear association between regions of extremely low recombination near to centromeres and high divergence in the data we present here. The factors that determine variability along chromosomes appear to be constant among distantly-related bee species (Jones et al. 2019). We therefore do not expect the landscape of recombination to be variable among the bumblebee species under investigation here, and expect the average rate to be similar to *B. terrestris*.

In addition to having low recombination rates, which enhances the effects of linked selection, IoDs in pericentromeric regions, as we observe here, could also be enhanced by genetic hitchhiking connected to the process of centromere drive (Crespi and Nosil 2013). This process results in rapid turnover of centromere satellite repeat sequences and proteins involved in the binding of centromeres to the spindle fibres during meiosis (Henikoff et al. 2001). Known speciation genes in *Drosophila*, *OdsH* (Bayes and Malik 2009) and *Zhr* (Sawamura et al. 1993), are both involved in interactions with satellite repeats. However, the IoDs are also gene rich, and therefore contain many functional sites that could be the focus of genetic hitchhiking or background selection, contributing to their elevated divergence.

Many studies of speciation genomics are focused on identifying specific genes that drive reproductive isolation by scanning the genome for IoDs (Ravinet et al. 2017). However, the

interaction between linked selection and recombination rate variation can explain the presence of IoDs in genomic comparisons. This has led some authors to distinguish between "incidental" IoDs formed by linked selection in regions of low recombination but irrelevant for the speciation process and true IoDs that harbour loci involved in reproductive isolation that are resistant to gene flow (Poelstra et al. 2014; Vijay et al. 2016). However, in this study we uncover evidence that IoDs that are found in regions of low recombination and appear to be resistant to gene flow in sympatry, indicating that the genome architecture is important in the formation of barriers to gene flow. The pervasiveness of this mechanism in nature is unclear. It is possible that a narrow set of conditions is required to generate IoDs with elevated d_{XY} indicating that there is often low statistical power to identify IoDs with elevated d_{XY} , particularly for species with short divergence times (Cruickshank and Hahn 2014). It is therefore possible that differential gene flow between the genomes of young species is more common than expected based on analysis of IoDs in pairwise genome comparisons. As there is no evidence that gene flow is ongoing between these species it was not possible to directly measure its effects across the genome. However, this study is consistent with a growing number of others indicating that selection against gene flow between incipient species can be highly polygenic, and strongly influenced by genome architecture (Michel et al. 2010; Coughlan and Matute 2020). It also supports a multitude of recent genome-wide studies that attest to the pervasiveness of gene flow and permeability of species barriers in nature.

Conclusion

We compared variation across the genomes of two recently-diverged cryptic bumblebee species living in sympatry. This comparison revealed the presence of restricted genomic islands (IoD) with elevated levels of absolute divergence (d_{XY}). This pattern suggests that the two species diverged under conditions of gene flow, which was restricted in regions of low recombination close to centromeres. These results imply that recombination rate variation could often be a crucial factor in determining the location of genomic barriers to gene flow between incipient species. We speculate that climatic fluctuations could be an important driver of speciation by this process in bumblebees with high-altitude habitats, whereby periods of warming lead to periodic population fragmentation at higher altitudes followed by secondary contact and differentiation under gene flow.

Methods

Genome sequencing and assembly

We generated a reference genome for *B. sylvicola* using ONT sequencing. DNA was extracted from a single male bee sampled from Niwot Ridge, CO, USA using a salt-isopropanol extraction followed by magnetic bead purification to remove fragments < 1000 bp and to concentrate the sample for library preparation. Sequencing was performed on a MinION with two R9.4 flowcells using the RAD004 kit (ONT) starting with 3-400 ng DNA per run, resulting in a yield of 9.4 Gbp with a total 2.5 million reads and a mean read length of 3.7 kbp. We used a multi-step approach to assemble the sequencing reads: downpore (<https://github.com/jteutenberg/downpore>) was used for adaptor trimming and splitting chimeric reads, trimmed reads were assembled using wtdbg2 using default settings (Ruan and Li 2020), then two rounds of the standalone consensus module Racon (<https://github.com/isovic/racon>) followed by further contig improvements with medaka v.0.4 (<https://github.com/nanoporetech/medaka>). For the medaka step, contigs of < 20 kbp were removed in order for the process to complete. The final polishing step involved two rounds of Pilon polishing (<https://github.com/broadinstitute/pilon>), whereby Illumina short reads were mapped to the assembly in order to correct the contigs around indels.

Long-range information from short-read sequencing of linked reads was obtained using 10x Genomics chromium technology. Sequencing was performed on the. A 10x GEM library was constructed from high-molecular-weight DNA from the same bee as for the ONP sequencing according to the manufacturer's recommended protocols. The resulting library was quantitated by qPCR and sequenced on one lane of a HiSeq 2500 using a HiSeq Rapid SBS sequencing kit version 2 to produce 150 bp paired-end sequences. We mapped the resultant reads to the assembly using Longranger v.2.1.4 and then ran Tigmint v1.1.2 to identify and correct errors in the assembly. ARCS+LINKS was used to scaffold the assembled contigs. We identified contigs that contained mitochondrial genes, and were therefore likely fragments of the mitochondrial genome, by running a blast search of *B. impatiens* mitochondrial genes across the assembly using BLAST+ v2.9.0. Any contigs containing two or more mitochondrial genes located within the expected distance of each other based on their locations on the mitochondrial genome were removed from the assembly, so that the final assembly did not contain partially assembled mitochondrial genome sequence. All contigs shorter than 10 kbp were also removed from the assembly. We ran BUSCO v3.0.2b (Simão et al. 2015) on the assembly in order to assess its completeness using the hymenoptera_odb9 lineage set and species *B. impatiens*. We performed whole-genome synteny alignments between the *B. terrestris* chromosome-level genome assembly and our *B. sylvicola* contigs using Satsuma v.3 (Grabherr et al. 2010) to arrange *B. sylvicola* contigs into pseudochromosomes, with the assumption of high structural conservation between the species. We performed both de novo and guided

transcriptome assemblies using reads from four different tissues: the abdomen, the head, the legs and the thorax. Full details of the annotation pipeline can be found in the Supplementary Methods.

Genome features

All genome features were calculated over 20 kbp non-overlapping windows for each contig. Genome GC content was measured using a custom perl script. We used GenMap (<https://github.com/cpockrandt/genmap>) to calculate mappability (uniqueness of k-mers) for each position in the genome, using a k-mer size (k) of 150 and a mismatch tolerance (e) of 2, and then averaged the output across windows using a custom perl script. RepeatMasker output from the genome annotation step was also summarised over windows, giving the proportion of each window that was characterised as repeat sequence, using a custom perl script.

In order to identify putative centromeric repeats in the *B. sylvicola* genome we ran centromere_seeker (https://github.com/cryanccampbell/centromere_seeker) on the raw ONT reads. This pipeline runs tandem repeats finder (TRF) (Benson 1999) to identify the longest and most prevalent tandem repeat arrays in the sequences, which are likely centromeric repeats. We used blastn (Altschul et al. 1990) to locate trimers of the identified 15 bp satellite in the genome to identify the likely locations of centromeres.

Population sampling

During the summer (July) of 2017 female worker bees from several species of *Pyrobombus* bumblebees were collected on seven mountains within the Rocky Mountains. Samples were collected with sweeping hand nets and kept in falcon tubes on cold packs in cool boxes for transport. Species identification was performed using a standard key (Williams et al. 2014). All samples were placed at -20°C for approximately 10 minutes before being dissected. Thoraces were stored in 95% ethanol for DNA extraction. Sampling effort achieved 217 *B. sylvicola*, 67 *B. incognitus*, and 17 *B. bifarius* individuals. Our sampling was supplemented with sequencing data for other *Pyrobombus* bees from published datasets available on the NCBI sequence read archive. These included 21 samples of *B. melanopygus* from western USA (Tian et al. 2019) (NCBI accession PRJNA526235), 4 extra samples of *B. bifarius* also collected in the Rocky Mountains, and 17 samples of *B. vancouverensis* from north-western USA (NCBI accession PRJNA592825).

Phenotypic variation among species

Measurements of intertegular distance (a proxy for body size) were made for all *B. sylvicola* and *B. incognitus* samples, as well as the 17 *B. bifarius* samples newly-collected in Colorado as part of this study (Supplementary Table S1). Intertegular distance was measured from scaled photographs of individual bees using ImageJ (<https://imagej.nih.gov/ij/>). The prementum and glossa were dissected from the head, mounted, and photographed for quantification using ImageJ. A subset of 69 samples defined as either *B. sylvicola* (N = 39) and *B. incognitus* (N = 30) by genetic clustering were randomly selected from across all sampling locations and examined morphologically in more detail. We characterised the shape of the malar space; the pile colour between the antenna and above the ocelli; the size, colour, and location of the ocelli relative to the supraorbital line; and the presence of black pile on the scutellum. We also recorded the colour of abdominal segments. Other body parts were sacrificed for genetic material; therefore, additional traits could not be characterized.

Population sequencing and variant calling

For all samples, DNA was extracted from the thorax of worker bees using the Qiagen Blood and Tissue kit. Paired-end sequencing libraries were prepared with Nextera Flex and samples were sequenced on an Illumina HiSeq X. Illumina paired end reads were mapped to our *B. sylvicola* genome assembly using the *mem* algorithm in BWA (Li and Durbin 2009). Mappings were piped to samtools (Li et al. 2009), where they were sorted by coordinate, written to bam files and indexed. Duplicate reads were marked and read groups were added in the bam files using the Picard suite of tools (<https://broadinstitute.github.io/picard/>). We used GATK to call variants, following their recommended Best Practices (<https://gatk.broadinstitute.org/hc/en-us>). Briefly, we ran HaplotypeCaller on each sample's bam file to create an individual-specific gVCF file. All gVCFs were then processed by GenomicsDBImport on a per contig basis, followed by GenotypeGVCFs to call variants. Variants were filtered using a hard set of filters using the VariantFiltration tool with thresholds recommended in the GATK best practices. These filters assess quality scores, depth of coverage, strand bias, and position on reads of variants to only retain high quality, high confidence SNPs. The resultant filtered vcf files were filtered for biallelic SNPs only. As all but one of the *B. vancouverensis* samples were haploid males, we used these samples to filter out SNPs that were called heterozygous in haploids and therefore are errors likely due to mapping issues at these sites.

Phylogenetic analysis of PEPCK and COI loci

In an attempt to identify the unknown species in our dataset, we extracted sequences of the PEPCK gene from our WGS datasets to generate per-sample sequences. We created a VCF file from all gVCFs containing only variants found across the PEPCK gene (located on contig_001: 5658366-5659292 bp)

and filtered these variants with the same set of filters detailed above. We then ran the GATK 'SelectVariants' tool to create a separate VCF per sample for the PEPCK gene, containing all positions where each sample differed from the reference. We then ran the GATK 'FastaAlternateReferenceMaker' tool on each of these VCFs to generate a sequence representing the PEPCK gene per sample based on the variants in each VCF. Manual inspection of mapping files was also performed to ensure sequences accurately reflected the evidence in the bam files. All sequences were concatenated into a single fasta file (see Supplementary Material), along with published PEPCK sequences for seven bumblebee species from (Martinet et al. 2019) and we generated a neighbour-joining tree from the sequences in SplitsTree4 v. 4.14.6. One thousand bootstrap replicates were performed to assess support.

As the mitochondrial genome was incompletely assembled in our genome assembly, we generated sequences from the COI locus using PCR and Sanger sequencing. We used the following primers adapted from (Hebert et al. 2004): LepF1, 5'-ATTCAACCAATCATAAAGATATTGG-3' and LepR1 5'-TAAACTTCTGGATGTCCAAAAATCA-3' for both PCR and sequencing. Thermal cycling conditions were as follows: one cycle of 1 min at 94°C, six cycles of 1 min at 94°C, 1 min and 30 sec at 45°C, and 1 min and 15 sec at 72°C, followed by 36 cycles of 1 min at 94°C, 1 min and 30 sec at 51°C, and 1 min and 15 sec at 72°C, with a final step of 5 min at 72°C. We processed sequences using CodonCode Aligner v. 9.0.1.3 and constructed phylogenetic trees in the same way as for the PEPCK locus.

Recombination rate variation

We measured patterns of linkage disequilibrium in our *B. sylvicola* population genomic dataset to infer the population-scale recombination rate, ρ , implemented in the program LDHat. A custom likelihood lookup table was created using the 'complete' program. The 'interval' program was used to estimate mean ρ across regions, running for 1.1 million iterations with chain sampling every 10,000 iterations, a burn-in of 100,000 iterations, and a block penalty of 1. Output from 'interval' was summarised using 'stat' and then converted into 20 kbp non-overlapping window averages across contigs using a custom perl script.

Genetic variation among species

We carried out a principal component analysis on a thinned set of likely independent SNPs. We retained one SNP every 10 kbp and a minimum minor allele frequency of 0.01 using vcfTools v. 0.1.15. The thinned vcf file was converted to a genlight object in the adegenet package in R v.4.0.2, via an intermediate plink raw format file. A PCA was carried out on the genlight object using the

glPca function of adegenet. A neighbour-joining tree was generated for the same set of SNPs in SplitsTree4 v. 4.14.6. We also generated a neighbor-net network (Bryant and Moulton 2004), implemented in SplitsTree4 v. 4.14.6, based on the same genome-wide SNP set in order to check whether the evolutionary history of *B. sylvicola* and *B. incognitus* appears tree-like and whether it presented any evidence for hybridisation.

Diversity and divergence

We calculated nucleotide diversity (π) within populations/species and relative (F_{ST}) and absolute (d_{XY}) divergence between populations/species in 20 kbp non-overlapping windows across the genome. We used scripts provided by Simon Martin (https://github.com/simonhmartin/genomics_general) for these analyses as they take into account the mixed ploidy among our samples. Per-window F_{ST} values were then standardised to a Z score (ZF_{ST}) in order to be able to compare genomic landscapes of divergence among pairs with different divergence times, using the following formula:

$$ZF_{ST} = \frac{\text{window } F_{ST} - \text{median } F_{ST}}{\text{standard deviation of } F_{ST}}$$

To measure branch-specific changes in allele frequency we calculated the population branch statistic (PBS) (Yi et al. 2010) for each species. For this we first obtained window-based F_{ST} measures for each species compared to an outgroup, *B. melanopygus*. All F_{ST} values for each 20 kbp window were then transformed to estimates of divergence times, T , using the equation (Nielsen et al. 1998):

$$T = -\ln(1 - F_{ST})$$

The length of the branch leading to each species could then be calculated using the following formula:

$$PBS = \frac{T_{S1S2} + T_{S1S0} - T_{S1S2}}{2}$$

Here, $S1$ and $S2$ refer to the two species being compared, and $S0$ is the outgroup *B. melanopygus*. The PBS value gives an estimate of the amount of divergence in terms of allele frequency change specific to a particular species ($S1$) since divergence from its common ancestor with the second species ($S2$).

Divergence time estimates

Under a model of neutral divergence of two populations from a common ancestor, F_{ST} can be converted into an estimate of time since divergence, T , where $T = t/3N_e$. Here, t is the number of generations since the two populations diverged and N_e is the effective size of each of the populations (Nielsen et al. 1998). Multiplying N_e by three is appropriate for haplodiploids. We estimated divergence time between *B. sylvicola* and *B. incognitus* and between *B. bifarius* and *B.*

vancouverensis by calculating T from F_{ST} across all regions of the genome sitting outside of the identified IoDs and then took the mean T . We estimated effective population sizes using an estimate of the population mutation rate, Watterson's estimator (θ_w), using the following equation:

$$\theta_w = \frac{K}{a_n}$$

Here, K is the number of segregating sites in the species and a_n is the $(n-1)^{\text{th}}$ harmonic number.

Values of θ_w were then used to calculate N_e for each species using the equation $3N_e = \theta/\mu$, where μ is the mutation rate. We used a value of $\mu = 3.6 \times 10^{-9}$, a direct estimate for *B. terrestris* (Liu et al. 2017). Multiplying T by $3N_e$ provided us with an estimate of the number of generations since the two species diverged (t). Assuming a generation time of one year, this estimate translates directly to the number of years since divergence. We calculated 95% confidence intervals around each estimate by bootstrapping the values of T from the 20 kbp window estimates using 5,000 bootstrap replicates with the *boot* package in R v.4.0.2.

Characterising islands of divergence

We characterised 20 kbp windows with ZF_{ST} values > 2 (2 standard deviations above the median) as highly divergent separately for each species. Highly divergent windows within 60 kbp of each other were then merged into single blocks. We classified divergent blocks greater than 100 kbp in length as IoDs in each pair. For the within-species and sympatric comparisons, any two IoDs within 1 Mbp of each other were merged into single IoDs as they likely are part of the same divergent region but small drops in ZF_{ST} in between meant they were not brought together in the previous step. We then defined all 20 kbp windows as either 'IoD' or 'background' for each population/species comparison separately and compared window measures of π , d_{XY} , PBS, recombination rate (ρ/kbp), GC content, mappability, and repeat content inside and outside of IoDs for each pair using Wilcoxon rank sum tests in R v.4.0.2.

We used permutation tests implemented with the R package *regioner* (Gel et al. 2016) to assess significance of overlap in the positions of IoDs between comparisons in a pairwise fashion. We used the 'randomizeRegions' function with the 'per.chromosome' option to randomize the location of each IoD along each pseudochromosome whilst maintaining its size. We performed 1,000 permutations and measured significance of the observed overlap by comparing it to a distribution of overlap in randomly positioned IoDs. Calculated Z-scores gave a measure of the strength of the association. We used this same method to assess whether exon content of IoDs is greater or less than that expected by chance, where the positions of exons and IoDs were randomised across the

genome whilst maintaining their size and 1,000 permutations were used to assess significance of the observed overlap.

To assess how diversity and divergence change when moving away from the centres of IoDs we took the positions of the centre of each IoD and used custom perl scripts to calculate average π , ZF_{ST} and d_{XY} in 20 kbp steps up to 5 Mbp away from each centre for the within-species, sympatric, and allopatric comparisons separately.

Measuring overlap between IoDs and centromere repeats

We measured overlap in the positions of IoDs and putative centromere repeat sequence for each comparison, using permutation tests in *regioneR* as described above. Here, the positions of both IoDs and repeats were randomised along each pseudochromosome whilst maintaining their size and 1,000 permutations were used to assess significance of the observed overlap. We used Z-scores to assess the effect size.

Uncovering evidence for structural variants between species

We used the program manta (Chen et al. 2016) to look for evidence of structural variants (SVs) between *B. sylvicola* and *B. incognitus* genomes. Manta uses mapped paired-end sequencing reads to discover, assemble and score large-scale SVs. In particular, we wanted to test whether there was evidence that the IoDs we identified could be the result of large inversions. As Manta is designed to run on a small set of samples, we ran it on the bam files of 12 *B. sylvicola* and 12 *B. incognitus* samples using default settings and then filtered the output for inversions that were fixed between the two species.

Data accessibility

The *Bombus sylvicola* genome assembly and Illumina whole genome resequencing reads for *B. sylvicola*, *B. incognitus*, and *B. bifarius* are available at NCBI under BioProject PRJNA646847. Alignments of the PEPCK locus generated from this data are available as Supplementary Information. COI sequences generated by this study are available at NCBI (accessions pending). All custom scripts used in the data analysis are available on GitHub at <https://github.com/MattChristmas/Pyrobombus-speciation>.

Acknowledgments

We thank Isabel Sullivan, Jacqueline Staab and Michael Osbourn who aided in specimen collection, identification and preservation and Lisa Danback, Kate Schleicher, Ying Thao, Anna Grobelny and Breana Cook who carried out phenotypic measurements. U.S. Forest Service provided access to field sites. We thank Baptiste Martinet for help with amplifying the COI locus. This research was funded by Swedish Research Council Formas grant 2016-00535 and a SciLifeLab National Biodiversity Project NP00046 to MTW. MK is financially supported by the Knut and Alice Wallenberg Foundation as part of the National Bioinformatics Infrastructure Sweden at SciLifeLab. The authors acknowledge support from the National Genomics Infrastructure in Stockholm and Uppsala funded by Science for Life Laboratory, the Knut and Alice Wallenberg Foundation and the Swedish Research Council. The computations and data handling were enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC) at UPPMAX partially funded by the Swedish Research Council through grant agreement no. 2018-05973.

References

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J. Mol. Biol.* 215:403–410.
- Bayes JJ, Malik HS. 2009. Altered heterochromatin binding by a hybrid sterility protein in *Drosophila* sibling species. *Science* 326:1538–1541.
- Begun DJ, Aquadro CF. 1992. Levels of naturally occurring DNA polymorphism correlate with recombination rates in *D. melanogaster*. *Nature* 356:519–520.
- Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* 27:573–580.
- Bertsch A, Schweer H, Titze A. 2004. Discrimination of the bumblebee species *Bombus lucorum*, *Bombus cryptarum* and *Bombus magnus* by morphological characters and male labial gland secretions (Hymenoptera: Apidae). *Beitr. Zur Entomol.* 54:365–386.
- Bossert S. 2015. Recognition and identification of species in the *Bombus lucorum*-complex—A review and outlook. *Dtsch. Entomol. Z.* 62:19–28.
- Brandvain Y, Kenney AM, Fligel L, Coop G, Sweigart AL. 2014. Speciation and introgression between *Mimulus nasutus* and *Mimulus guttatus*. *PLOS Genet.* 10:e1004410.
- Bryant D, Moulton V. 2004. Neighbor-Net: An Agglomerative Method for the Construction of Phylogenetic Networks. *Mol. Biol. Evol.* 21:255–265.
- Burri R, Nater A, Kawakami T, Mugal CF, Olason PI, Smeds L, Suh A, Dutoit L, Bureš S, Garamszegi LZ, et al. 2015. Linked selection and recombination rate variation drive the evolution of the genomic landscape of differentiation across the speciation continuum of *Ficedula* flycatchers. *Genome Res.* 25:1656–1665.

- Cameron SA, Hines HM, Williams PH. 2007. A comprehensive phylogeny of the bumble bees (*Bombus*). *Biol. J. Linn. Soc.* 91:161–188.
- Chapman MA, Hiscock SJ, Filatov DA. 2016. The genomic bases of morphological divergence and reproductive isolation driven by ecological speciation in *Senecio* (Asteraceae). *J. Evol. Biol.* 29:98–113.
- Charlesworth B, Nordborg M, Charlesworth D. 1997. The effects of local selection, balanced polymorphism and background selection on equilibrium patterns of genetic diversity in subdivided populations. *Genet. Res.* 70:155–174.
- Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Källberg M, Cox AJ, Kruglyak S, Saunders CT. 2016. Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* 32:1220–1222.
- Coughlan JM, Matute DR. 2020. The importance of intrinsic postzygotic barriers throughout the speciation process. *Philos. Trans. R. Soc. B Biol. Sci.* 375:20190533.
- Coyne JA, Orr HA. 2004. Speciation. Sinauer
- Crespi B, Nosil P. 2013. Conflictual speciation: species formation via genomic conflict. *Trends Ecol. Evol.* 28:48–57.
- Cruickshank TE, Hahn MW. 2014. Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow. *Mol. Ecol.* 23:3133–3157.
- Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Künstner A, Mäkinen H, Nadachowska-Brzyska K, Qvarnström A, et al. 2012. The genomic landscape of species divergence in *Ficedula* flycatchers. *Nature* 491:756–760.
- Ezray BD, Wham DC, Hill C, Hines HM. 2019. Müllerian mimicry in bumble bees is a transient continuum. *bioRxiv*:513275.
- Feder JL, Egan SP, Nosil P. 2012. The genomics of speciation-with-gene-flow. *Trends Genet.* 28:342–350.
- Feulner PGD, Chain FJJ, Panchal M, Huang Y, Eizaguirre C, Kalbe M, Lenz TL, Samonte IE, Stoll M, Bornberg-Bauer E, et al. 2015. Genomics of divergence along a continuum of parapatric population differentiation. *PLOS Genet.* 11:e1004966.
- Gel B, Díez-Villanueva A, Serra E, Buschbeck M, Peinado MA, Malinverni R. 2016. regioneR: an R/Bioconductor package for the association analysis of genomic regions based on permutation tests. *Bioinformatics* 32:289–291.
- Ghisbain G, Lozier JD, Rahman SR, Ezray BD, Tian L, Ulmer JM, Heraghty SD, Strange JP, Rasmont P, Hines HM. 2020. Substantial genetic divergence and lack of recent gene flow support cryptic speciation in a colour polymorphic bumble bee (*Bombus bifarius*) species complex. *Syst. Entomol.* 45:635–652.
- Goulson D. 2003. Bumblebees: their behaviour and ecology. Oxford University Press, USA

- Grabherr MG, Russell P, Meyer M, Mauceli E, Alföldi J, Di Palma F, Lindblad-Toh K. 2010. Genome-wide synteny through highly sensitive sequence alignment: Satsuma. *Bioinformatics* 26:1145–1151.
- Guerrero RF, Hahn MW. 2017. Speciation as a sieve for ancestral polymorphism. *Mol. Ecol.* 26:5362–5368.
- Han F, Lamichhaney S, Grant BR, Grant PR, Andersson L, Webster MT. 2017. Gene flow, ancient polymorphism, and ecological adaptation shape the genomic landscape of divergence among Darwin’s finches. *Genome Res.* 27:1004–1015.
- Hebert PDN, Penton EH, Burns JM, Janzen DH, Hallwachs W. 2004. Ten species in one: DNA barcoding reveals cryptic species in the neotropical skipper butterfly *Astraptes fulgerator*. *Proc. Natl. Acad. Sci.* 101:14812–14817.
- Henikoff S, Ahmad K, Malik HS. 2001. The centromere paradox: stable inheritance with rapidly evolving DNA. *Science* 293:1098–1102.
- Hewitt G. 2000. The genetic legacy of the Quaternary ice ages. *Nature* 405:907–913.
- Hines HM. 2008. Historical biogeography, divergence times, and diversification patterns of bumble bees (Hymenoptera: Apidae: Bombus). *Syst. Biol.* 57:58–75.
- Hines HM, Cameron SA, Williams PH. 2006. Molecular phylogeny of the bumble bee subgenus *Pyrobombus* (Hymenoptera : Apidae : Bombus) with insights into gene utility for lower-level analysis. *Invertebr. Syst.* 20:289–303.
- Irwin DE, Milá B, Toews DPL, Brelsford A, Kenyon HL, Porter AN, Grossen C, Delmore KE, Alcaide M, Irwin JH. 2018. A comparison of genomic islands of differentiation across three young avian species pairs. *Mol. Ecol.* 27:4839–4855.
- Janoušek V, Munclinger P, Wang L, Teeter KC, Tucker PK. 2015. Functional organization of the genome may shape the species boundary in the house mouse. *Mol. Biol. Evol.* 32:1208–1220.
- Jones JC, Wallberg A, Christmas MJ, Kapheim KM, Webster MT. 2019. Extreme differences in recombination rate between the genomes of a solitary and a social bee. *Mol. Biol. Evol.* 36:2277–2291.
- Juric I, Aeschbacher S, Coop G. 2016. The strength of selection against neanderthal introgression. *PLOS Genet.* 12:e1006340.
- Kawakami T, Wallberg A, Olsson A, Wintermantel D, Miranda JR de, Allsopp M, Rundlöf M, Webster MT. 2019. Substantial heritable variation in recombination rate on multiple scales in honeybees and bumblebees. *Genetics* 212:1101–1119.
- Kirkpatrick M, Ravigné V. 2002. Speciation by Natural and Sexual Selection: Models and Experiments. *Am. Nat.* 159:S22–S35.
- Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* 37:540–546.

- Lamichhane S, Berglund J, Almén MS, Maqbool K, Grabherr M, Martinez-Barrio A, Promerová M, Rubin C-J, Wang C, Zamani N, et al. 2015. Evolution of Darwin's finches and their beaks revealed by genome sequencing. *Nature* 518:371–375.
- Lecocq T, Dellicour S, Michez D, Lhomme P, Vanderplanck M, Valterová I, Rasplus J-Y, Rasmont P. 2013. Scent of a break-up: phylogeography and reproductive trait divergences in the red-tailed bumblebee (*Bombus lapidarius*). *BMC Evol. Biol.* 13:263.
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* 25:1754–1760.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079.
- Liu H, Jia Y, Sun X, Tian D, Hurst LD, Yang S. 2017. Direct Determination of the Mutation Rate in the Bumblebee Reveals Evidence for Weak Recombination-Associated Mutation and an Approximate Rate Constancy in Insects. *Mol. Biol. Evol.* 34:119–130.
- Liu X, Glémin S, Karrenberg S. 2020. Evolution of putative barrier loci at an intermediate stage of speciation with gene flow in campions (*Silene*). *Mol. Ecol.* 29:3511–3525.
- Malinsky M, Challis RJ, Tyers AM, Schiffels S, Terai Y, Ngatunga BP, Miska EA, Durbin R, Genner MJ, Turner GF. 2015. Genomic islands of speciation separate cichlid ecomorphs in an East African crater lake. *Science* 350:1493–1498.
- Mallet J, Besansky N, Hahn MW. 2016. How reticulated are species? *BioEssays* 38:140–149.
- Martin SH, Dasmahapatra KK, Nadeau NJ, Salazar C, Walters JR, Simpson F, Blaxter M, Manica A, Mallet J, Jiggins CD. 2013. Genome-wide evidence for speciation with gene flow in *Heliconius* butterflies. *Genome Res.* 23:1817–1828.
- Martin SH, Davey JW, Salazar C, Jiggins CD. 2019. Recombination rate variation shapes barriers to introgression across butterfly genomes. *PLOS Biol.* 17:e2006288.
- Martinet B, Lecocq T, Brasero N, Gerard M, Urbanová K, Valterová I, Gjershaug JO, Michez D, Rasmont P. 2019. Integrative taxonomy of an arctic bumblebee species complex highlights a new cryptic species (*Apidae: Bombus*). *Zool. J. Linn. Soc.* 187:599–621.
- McVean GAT, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. 2004. The fine-scale structure of recombination rate variation in the human genome. *Science* 304:581–584.
- Michel AP, Sim S, Powell THQ, Taylor MS, Nosil P, Feder JL. 2010. Widespread genomic divergence during sympatric speciation. *Proc. Natl. Acad. Sci.* 107:9724–9729.
- Murray TE, Fitzpatrick Ú, Brown MJF, Paxton RJ. 2008. Cryptic species diversity in a widespread bumble bee complex revealed using mitochondrial DNA RFLPs. *Conserv. Genet.* 9:653–666.
- Nachman MW, Payseur BA. 2012. Recombination rate variation and speciation: theoretical predictions and empirical results from rabbits and mice. *Philos. Trans. R. Soc. B Biol. Sci.* 367:409–421.

- Nielsen R, Mountain JL, Huelsenbeck JP, Slatkin M. 1998. Maximum-Likelihood Estimation of Population Divergence Times and Population Phylogeny in Models without Mutation. *Evolution* 52:669–677.
- Nordborg M, Charlesworth B, Charlesworth D. 1996. The effect of recombination on background selection. *Genet. Res.* 67:159–174.
- Papadopulos AST, Igea J, Dunning LT, Osborne OG, Quan X, Pellicer J, Turnbull C, Hutton I, Baker WJ, Butlin RK, et al. 2019. Ecological speciation in sympatric palms: 3. Genetic map reveals genomic islands underlying species divergence in *Howea*. *Evolution* 73:1986–1995.
- Pessia E, Popa A, Mousset S, Rezvoy C, Duret L, Marais GAB. 2012. Evidence for widespread GC-biased gene conversion in eukaryotes. *Genome Biol. Evol.* 4:675–682.
- Poelstra JW, Vijay N, Bossu CM, Lantz H, Ryll B, Müller I, Baglione V, Unneberg P, Wikelski M, Grabherr MG, et al. 2014. The genomic landscape underlying phenotypic integrity in the face of gene flow in crows. *Science* 344:1410–1414.
- Powell DL, García-Olazábal M, Keegan M, Reilly P, Du K, Díaz-Loyo AP, Banerjee S, Blakkan D, Reich D, Andolfatto P, et al. 2020. Natural hybridization reveals incompatible alleles that cause melanoma in swordtail fish. *Science* 368:731–736.
- Ravinet M, Faria R, Butlin RK, Galindo J, Bierne N, Rafajlović M, Noor M a. F, Mehlig B, Westram AM. 2017. Interpreting the genomic landscape of speciation: a road map for finding barriers to gene flow. *J. Evol. Biol.* 30:1450–1477.
- Renaut S, Grassa CJ, Yeaman S, Moyers BT, Lai Z, Kane NC, Bowers JE, Burke JM, Rieseberg LH. 2013. Genomic islands of divergence are not affected by geography of speciation in sunflowers. *Nat. Commun.* 4:1827.
- Roesti M, Hendry AP, Salzburger W, Berner D. 2012. Genome divergence during evolutionary diversification as revealed in replicate lake–stream stickleback population pairs. *Mol. Ecol.* 21:2852–2862.
- Ruan J, Li H. 2020. Fast and accurate long-read assembly with wtdbg2. *Nat. Methods* 17:155–158.
- Rundle HD, Nosil P. 2005. Ecological speciation. *Ecol. Lett.* 8:336–352.
- Sawamura K, Yamamoto MT, Watanabe TK. 1993. Hybrid lethal systems in the *Drosophila melanogaster* species complex. II. The Zygotic hybrid rescue (*Zhr*) gene of *D. melanogaster*. *Genetics* 133:307–313.
- Schumer M, Cui R, Powell DL, Dresner R, Rosenthal GG, Andolfatto P. 2014. High-resolution mapping reveals hundreds of genetic incompatibilities in hybridizing fish species. McVean G, editor. *eLife* 3:e02535.
- Schumer M, Xu C, Powell DL, Durvasula A, Skov L, Holland C, Blazier JC, Sankararaman S, Andolfatto P, Rosenthal GG, et al. 2018. Natural selection interacts with recombination to shape the evolution of hybrid genomes. *Science* 360:656–660.
- Seehausen O, Butlin RK, Keller I, Wagner CE, Boughman JW, Hohenlohe PA, Peichel CL, Saetre G-P, Bank C, Brännström Å, et al. 2014. Genomics and the origin of species. *Nat. Rev. Genet.* 15:176–192.

- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31:3210–3212.
- Soria-Carrasco V, Gompert Z, Comeault AA, Farkas TE, Parchman TL, Johnston JS, Buerkle CA, Feder JL, Bast J, Schwander T, et al. 2014. Stick insect genomes reveal natural selection's role in parallel speciation. *Science* 344:738–742.
- Stankowski S, Chase MA, Fuiten AM, Rodrigues MF, Ralph PL, Streisfeld MA. 2019. Widespread selection and gene flow shape the genomic landscape during a radiation of monkeyflowers. *PLOS Biol.* 17:e3000391.
- Talla V, Kalsoom F, Shipilina D, Marova I, Backström N. 2017. Heterogeneous patterns of genetic diversity and differentiation in european and siberian chiffchaff (*Phylloscopus collybita abietinus*/P. tristis). *G3 Genes Genomes Genet.* 7:3983–3998.
- Tian L, Rahman SR, Ezray BD, Franzini L, Strange JP, Lhomme P, Hines HM. 2019. A homeotic shift late in development drives mimetic color variation in a bumble bee. *Proc. Natl. Acad. Sci.*:201900365.
- Turner TL, Hahn MW. 2010. Genomic islands of speciation or genomic islands and speciation? *Mol. Ecol.* 19:848–850.
- Turner TL, Hahn MW, Nuzhdin SV. 2005. Genomic islands of speciation in *Anopheles gambiae*. *PLOS Biol.* 3:e285.
- Uemura R, Motoyama H, Masson-Delmotte V, Jouzel J, Kawamura K, Goto-Azuma K, Fujita S, Kuramoto T, Hirabayashi M, Miyake T, et al. 2018. Asynchrony between Antarctic temperature and CO₂ associated with obliquity over the past 720,000 years. *Nat. Commun.* 9:961.
- Via S. 2012. Divergence hitchhiking and the spread of genomic isolation during ecological speciation-with-gene-flow. *Philos. Trans. R. Soc. B Biol. Sci.* 367:451–460.
- Via S, West J. 2008. The genetic mosaic suggests a new role for hitchhiking in ecological speciation. *Mol. Ecol.* 17:4334–4345.
- Vijay N, Bossu CM, Poelstra JW, Weissensteiner MH, Suh A, Kryukov AP, Wolf JBW. 2016. Evolution of heterogeneous genome differentiation across multiple contact zones in a crow species complex. *Nat. Commun.* 7:13195.
- Vimeux F, Cuffey KM, Jouzel J. 2002. New insights into southern hemisphere temperature changes from Vostok ice cores using deuterium excess correction. *Earth Planet. Sci. Lett.* 203:829–843.
- Wallberg A, Han F, Wellhagen G, Dahle B, Kawata M, Haddad N, Simões ZLP, Allsopp MH, Kandemir I, De la Rúa P, et al. 2014. A worldwide survey of genome sequence variation provides insight into the evolutionary history of the honeybee *Apis mellifera*. *Nat. Genet.* 46:1081–1088.
- Wilfert L, Gadau J, Schmid-Hempel P. 2007. Variation in genomic recombination rates among animal taxa and the case of social insects. *Heredity* 98:189–197.

Williams P. 2007. The distribution of bumblebee colour patterns worldwide: possible significance for thermoregulation, crypsis, and warning mimicry. *Biol. J. Linn. Soc.* 92:97–118.

Williams PH, Thorp RW, Richardson LL, Colla S. 2014. Bumble bees of North America: an identification guide. Princeton University Press

Wu C-I. 2001. The genic view of the process of speciation. *J. Evol. Biol.* 14:851–865.

Yeaman S. 2013. Genomic rearrangements and the evolution of clusters of locally adaptive loci. *Proc. Natl. Acad. Sci.* 110:E1743–E1751.

Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, et al. 2010. Sequencing of 50 Human Exomes Reveals Adaptation to High Altitude. *Science* 329:75–78.

Tables

Table 1. Summary statistics for the four *Pyrobombus* species

Species	N	Number of chromosomes	Number of SNPs	Effective population size, N_e	Watterson's theta per base (θ_w)	Nucleotide diversity (π)
<i>B. sylvicola</i>	217	434	4,655,117	260,000	0.0028	0.0025
<i>B. incognitus</i>	67	134	4,891,459	320,000	0.0035	0.0028
<i>B. bifarius</i>	21	41	1,924,407	164,000	0.0018	0.0014
<i>B. vancouverensis</i>	17	18	2,057,379	216,000	0.0023	0.0016

Table 2. Differences in genomic content inside and outside of Islands of Divergence (IoDs)

Metric	Recombination rate (ρ /kbp)		GC content (%)		Repeat content (%)		Mappability		Exonic sequence (%)	
	In	Out	In	Out	In	Out	In	Out	In	Out
IoDs										
Within-species	2.3***	21.7	36.6***	38.1	21.0***	12.2	0.97***	0.99	13.4	15.7
Sympatric	2.7***	23.8	36.4***	38.3	16.1***	12.0	0.97***	0.99	15.8	15.5
Allopatric	11.1***	21.5	34.4***	38.3	18.4***	12.1	0.97***	0.99	29.0***	14.6

*** $p < 0.001$, ** $p < 0.01$. Significance of difference in recombination rate, proportion GC content, repeat content, and mappability inside and outside of IoDs assessed with Wilcoxon rank sum test. Significance of proportion of exonic sequence inside of IoDs assessed using permutation tests.

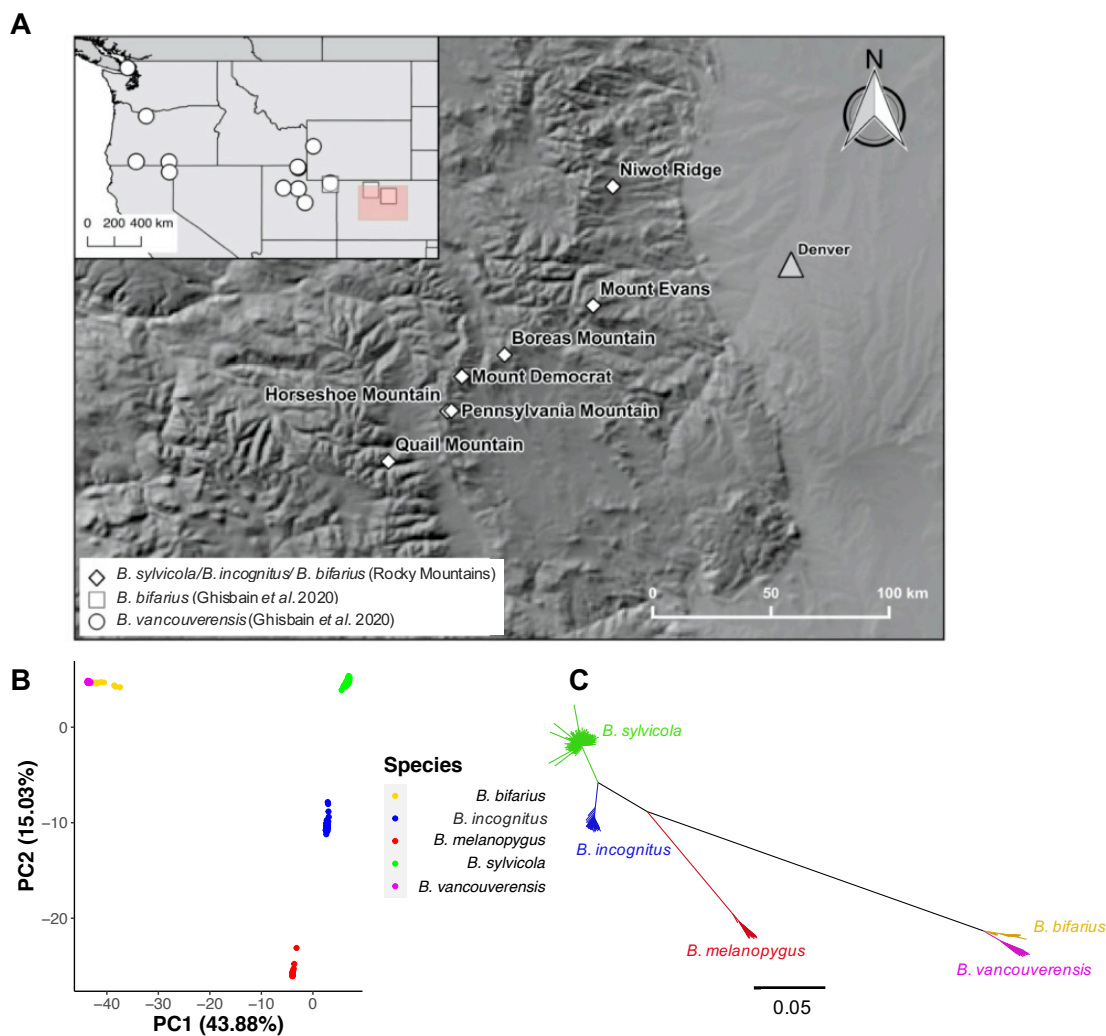


Figure 1. Sampling locations and genomic variation in *Pyrobombus* bumblebees. (A) Map showing the seven locations in Colorado where *Pyrobombus* bees were sampled for this study, as well as the sampling locations of *B. bifarius* and *B. vancouverensis* from a previous study (inset) (Ghisbain et al. 2020). *B. melanopygus* was collected widely across western USA in a previous study (Tian et al. 2019). (B) Principal component analysis and (C) a neighbour-joining tree based on genome-wide SNPs thinned for one SNP every 10 kbp of the five *Pyrobombus* species included in this study revealed distinct genetic divergence between *B. sylvicola* and *B. incognitus*. Scale bar on tree represents sequence divergence (%).

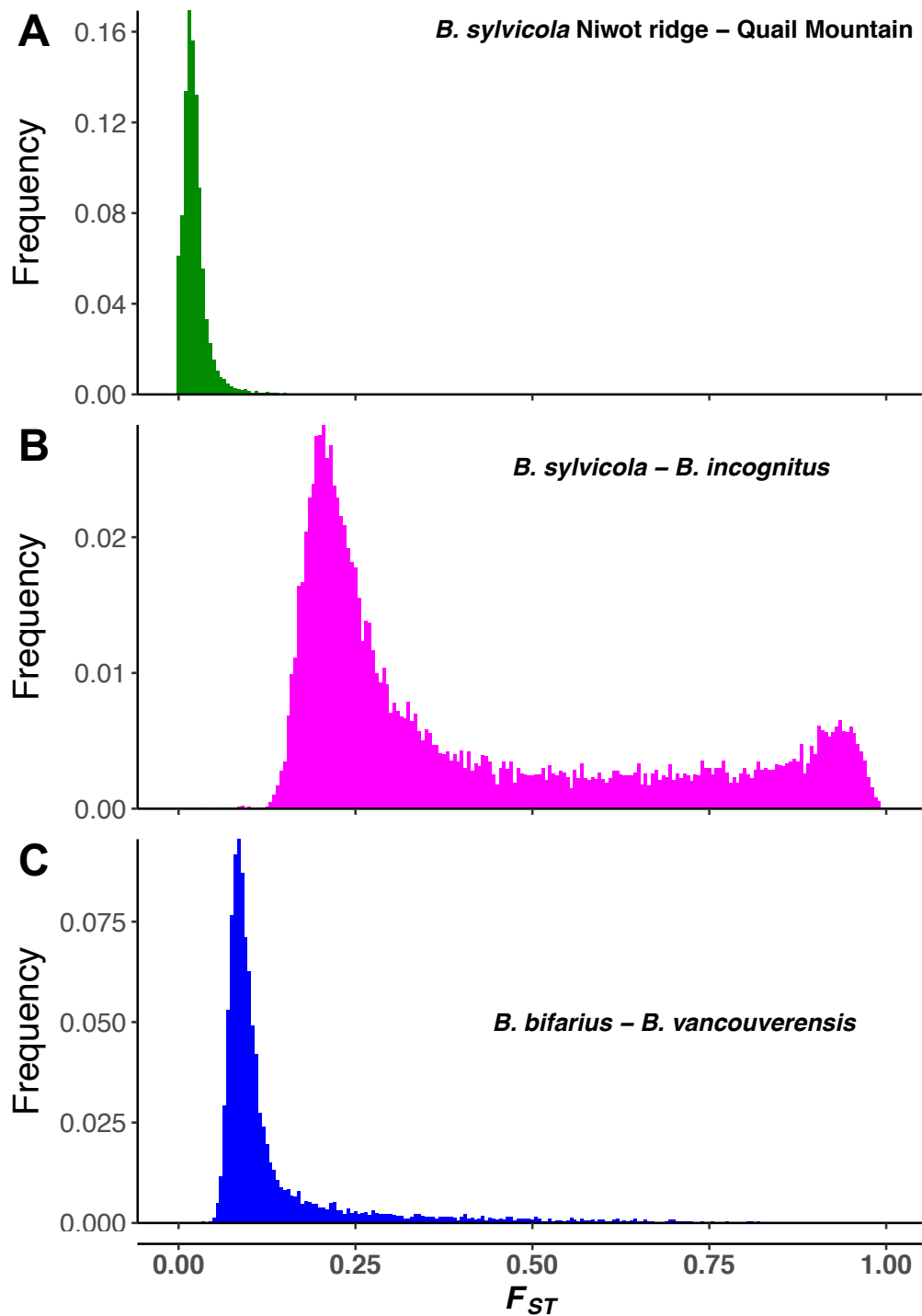


Figure 2. Histograms showing distributions of F_{ST} measured in 20 kbp windows across the genome for (A) a within-species comparison of two *Bombus sylvicola* populations, Niwot Ridge (n=43) and Quail Mountain (n=17), (B) a sympatric comparison of *B. sylvicola* (n=217), and *B. incognitus* (n=67), and (C) an allopatric comparison of *B. bifarius* (n=21), and *B. vancouverensis* (n=17). Note the distinct bimodal distribution of the sympatric comparison, indicating differential divergence across the genome.

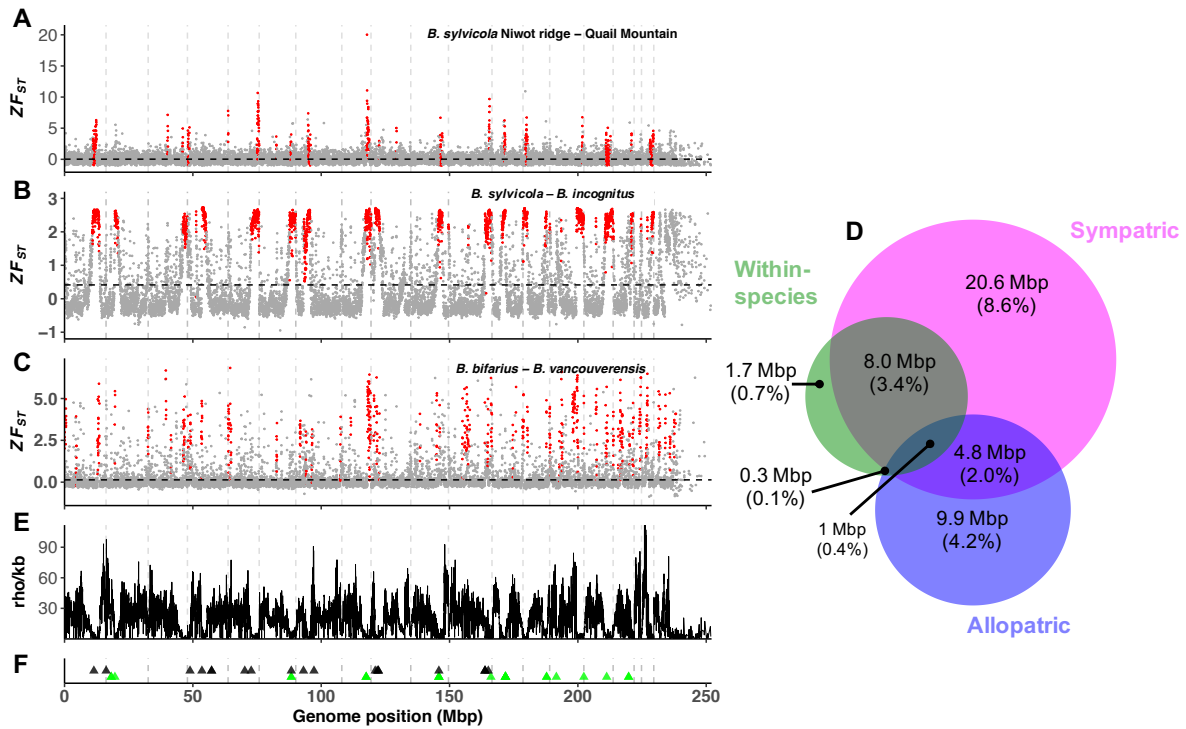


Figure 3. Genome-wide Z_{FST} scores measured in 20 kbp non-overlapping sliding windows for (A) within-species *B. sylvicola* Niwot Ridge – Quail Mountain comparison, (B) sympatric *B. sylvicola* – *B. incognitus* comparison, and (C) allopatric *B. bifarius* – *B. vancouverensis* comparison. Red dots represent 20 kbp windows that are located in islands of divergence (IoDs), defined as regions of extreme divergence (Z_{FST} scores > 2) over 100 kbp in length. Horizontal black dashed lines represent mean values. (D) Venn diagram showing overlap in positions of IoDs between comparisons, including size of overlap in Mbp and as a percentage of the genome. (E) Recombination rate variation across the genome measured in ρ /kbp. (F) Positions of putative centromere tandem repeat arrays, where black triangles indicate arrays > 1 kbp and green indicates arrays < 1 kbp. Vertical grey dashed lines represent boundaries of the 18 pseudo-chromosomes, with data from unplaced contigs shown to the right of the last pseudo-chromosome.

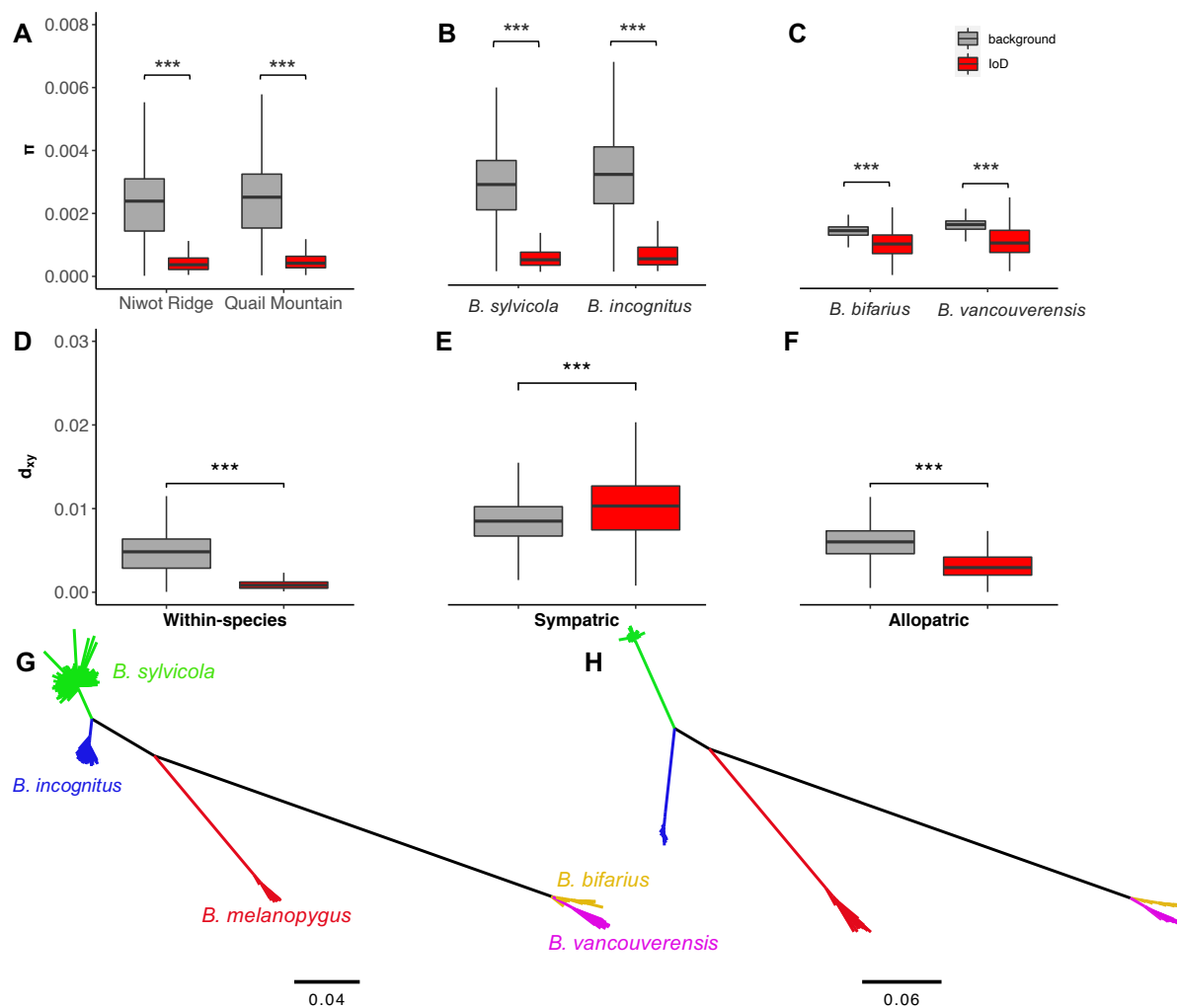


Figure 4. Differences in nucleotide diversity (π) and absolute divergence (d_{xy}) inside and outside of islands of divergence (IoDs, where $ZF_{ST} > 2$) for (A,D) a within-species comparison of two *B. sylvicola* populations (Niwot Ridge and Quail Mountain), (B,E) a sympatric comparison of *B. sylvicola* and *B. incognitus*, and (C,F) an allopatric comparison of *B. bifarius* and *B. vancouverensis*. Neighbour-joining trees based on single nucleotide polymorphisms (G) outside of IoDs and (H) inside IoDs in the sympatric comparison. Scale bars on trees represents sequence divergence (%).

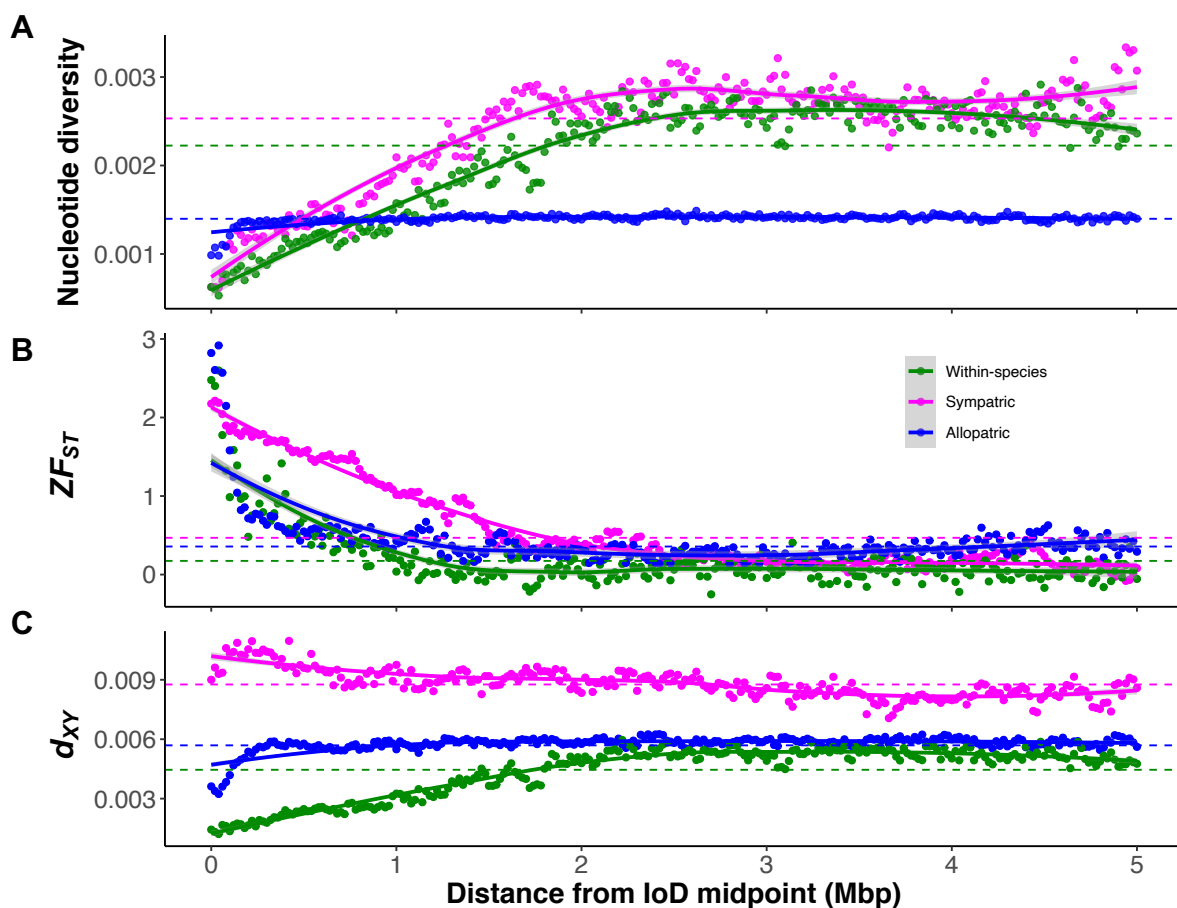


Figure 5. Average changes in (A) nucleotide diversity (π), (B) relative divergence (ZF_{ST}), and (C) absolute divergence (d_{xy}) stepping away from centres of Islands of Divergence in 20 kbp windows in both directions. Within-species comparisons are Niwot ridge and Quail Mountain populations of *Bombus sylvicola*, the sympatric comparison is *B. sylvicola* and *B. incognitus*, and the allopatric comparison is *B. bifarius* and *B. vancouverensis*. For nucleotide diversity, data for only one species from each pair is shown for clarity: Within-species = Niwot Ridge, sympatric = *B. sylvicola*, allopatric = *B. bifarius*. Dashed lines represent mean values. Smooth curves are based on locally estimated scatterplot smoothing (LOESS), with 95% confidence intervals shown in grey.